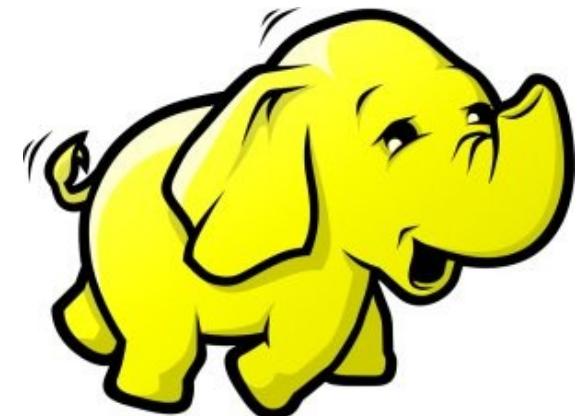




# 淺談海量資料的趨勢、挑戰與因應對策

Big Data : the Trends, Challenges and Solutions

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Agenda 演講大綱

**What is Big Data ?**

何謂海量資料

**Why should we care?**

爲何需要關切

**When to deploy it ?**

何時導入技術

**How to handle it ?**

三大因應策略

**Who is key player ?**

誰是成功關鍵

# WHAT



## What is Big Data ?

## 何謂海量資料

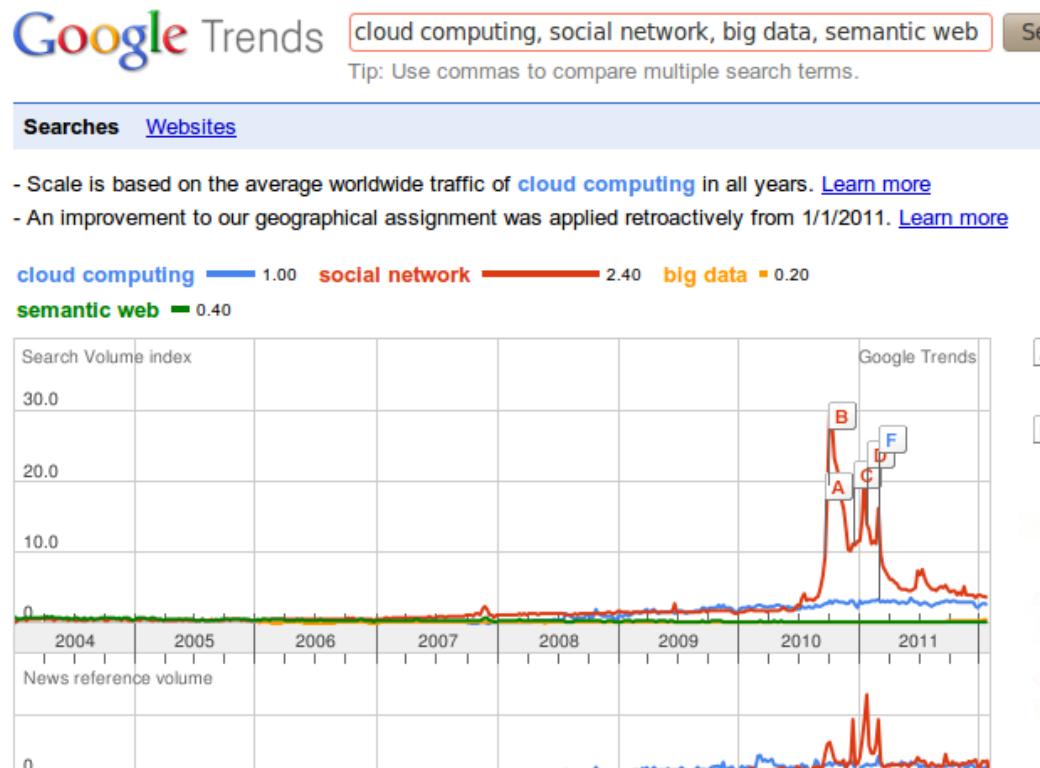
趨勢  
Trends

定義  
Definitions

挑戰：管理維度  
The Six Dimensions

Source: <http://www.2010taipeiexpo.tw/ct.asp?xItem=17186&CtNode=5952&mp=3>

# Trends .... It's all about **Buzzwords** .... 「趨勢」亦或「流行語」？ Web 3.0, Cloud Computing, Social Network, Big Data, ....



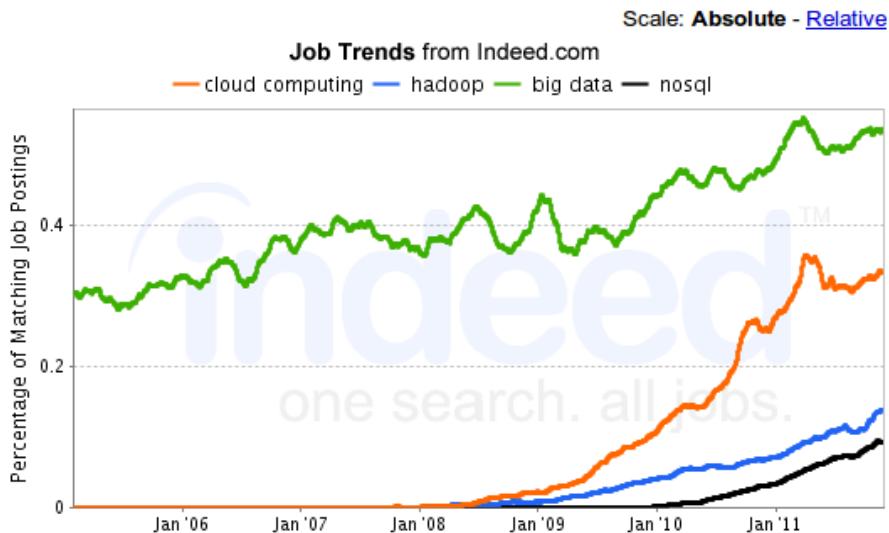
語意網（ Semantic Web ）從 2001 年開始制定標準後，逐漸下滑。而同義詞 Web 3.0 也呈現相似趨勢。海量資料（ Big Data ）與其關鍵技術 Hadoop ，則仍在上揚中。



整體而言，雲端運算（ Cloud Computing ）與社交網路（ Social Network ）呈現上揚。且社交網路比雲端運算還引人注目。

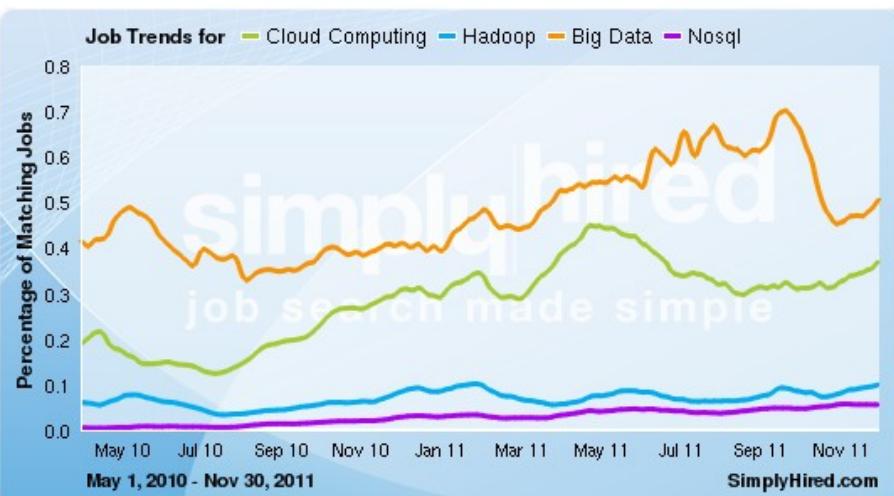
# Trends of Market Needs 市場需求趨勢

cloud computing, hadoop, big data, nosql Job Trends



Indeed.com searches millions of jobs from thousands of job sites.  
This job trends graph shows the percentage of jobs we find that contain your search terms.

Find [Cloud Computing jobs](#), [Hadoop jobs](#), [Big Data jobs](#), [Nosql jobs](#)



美國軟體就業市場分析，根據 indeed 與 simply hired 兩間公司的趨勢觀察，都得到一樣的結果：  
**Big Data > Cloud Computing > Hadoop > NoSQL**

CIO technologies	Ranking of technologies CIOs selected as one of their top 3 priorities in 2012			
	2012	2011	2010	2009
Analytics and business intelligence	1	5	5	1
Mobile technologies	2	3	6	12
Cloud computing (SaaS, IaaS, PaaS)	3	1	2	16
Collaboration technologies (workflow)	4	8	11	5
Virtualization	5	2	1	3
Legacy modernization	6	7	15	4
IT management	7	4	10	*
Customer relationship management	8	18	*	*
ERP applications	9	13	14	2
Security	10	12	9	8
Social media/Web 2.0	11	10	3	15

Gartner CIO Agenda 2012 前三名：

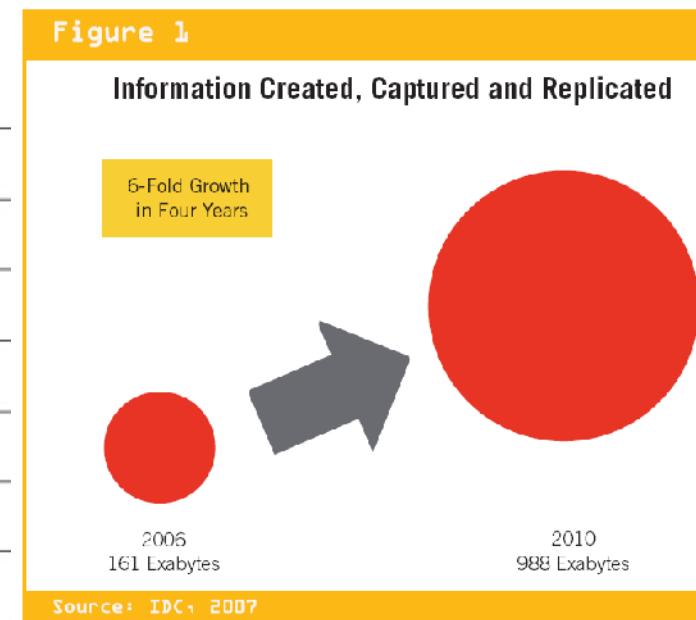
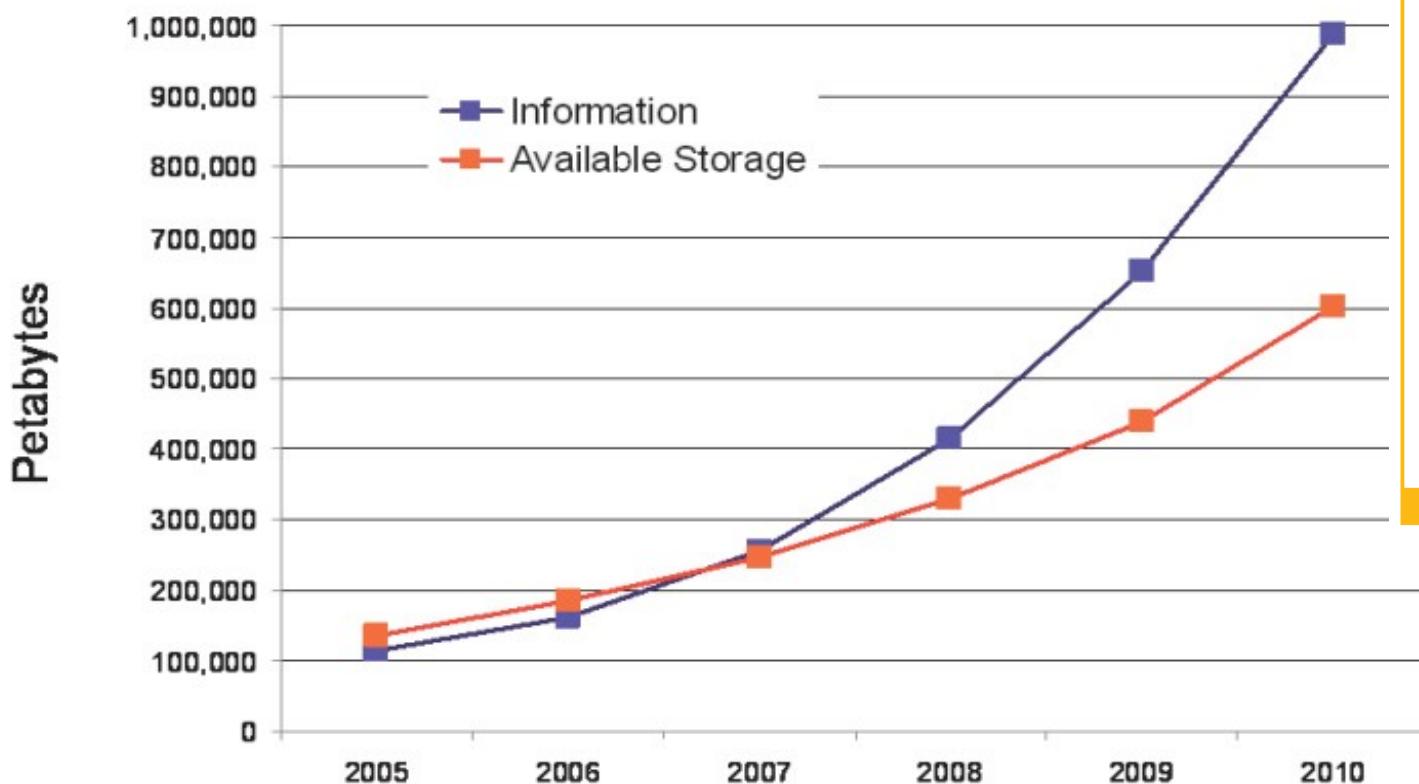
- [1] Business Intelligence (Big Data)
- [2] Mobile technology
- [3] Cloud Computing

# How BIG? 讓我們先來認識一下容量單位

Bit (b)	1 or 0
Byte (B)	8 bits
Kilobyte (KB)	1,000 bytes
Megabyte (MB)	1,000 KB
Gigabyte (GB)	1,000 MB
Terabyte (TB)	1,000, GB
Petabyte (PB)	1,000 TB
Exabyte (EB)	1,000 PB
Zettabyte (ZB)	1,000 EB

# Data Explosion!! 始於 2007 的「資料大爆炸」時代

## Information Versus Available Storage



2007 年，IDC 預估  
2010 年會成長六倍！  
(相較 2006 年)

Source: IDC, 2007

出處：[The Expanding Digital Universe](#),

A Forecast of Worldwide Information Growth Through 2010,  
March 2007, An IDC White Paper - sponsored by EMC

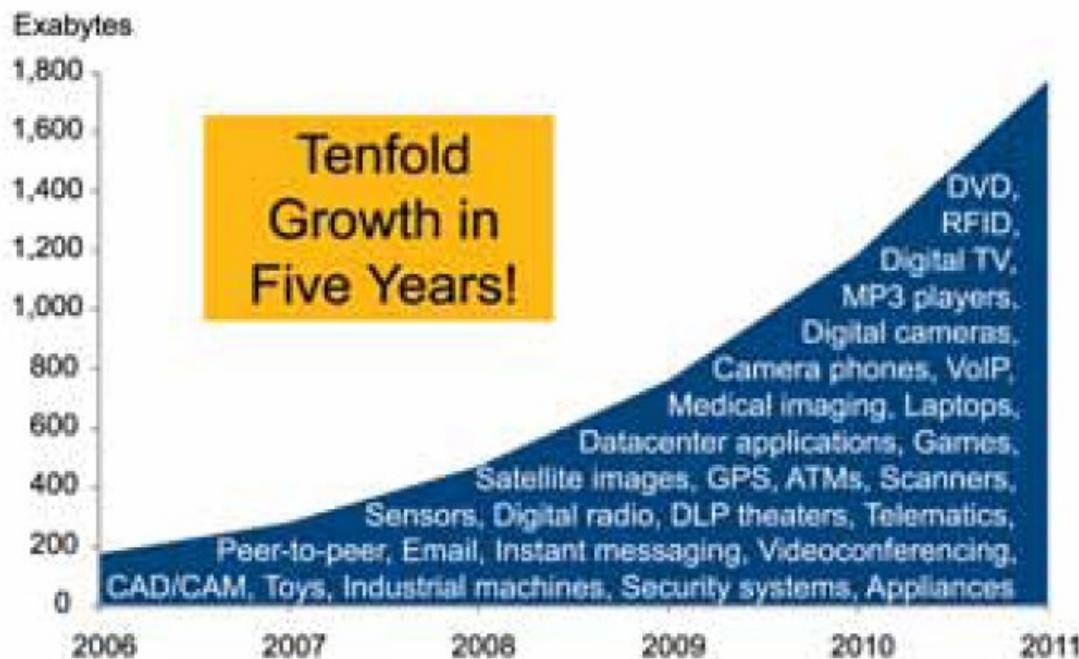
<http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>

2006 161 EB  
2010 988 EB (預測)

# Data Explosion!! 始於 2007 的「資料大爆炸」時代

Figure 1

Digital Information Created, Captured, Replicated Worldwide



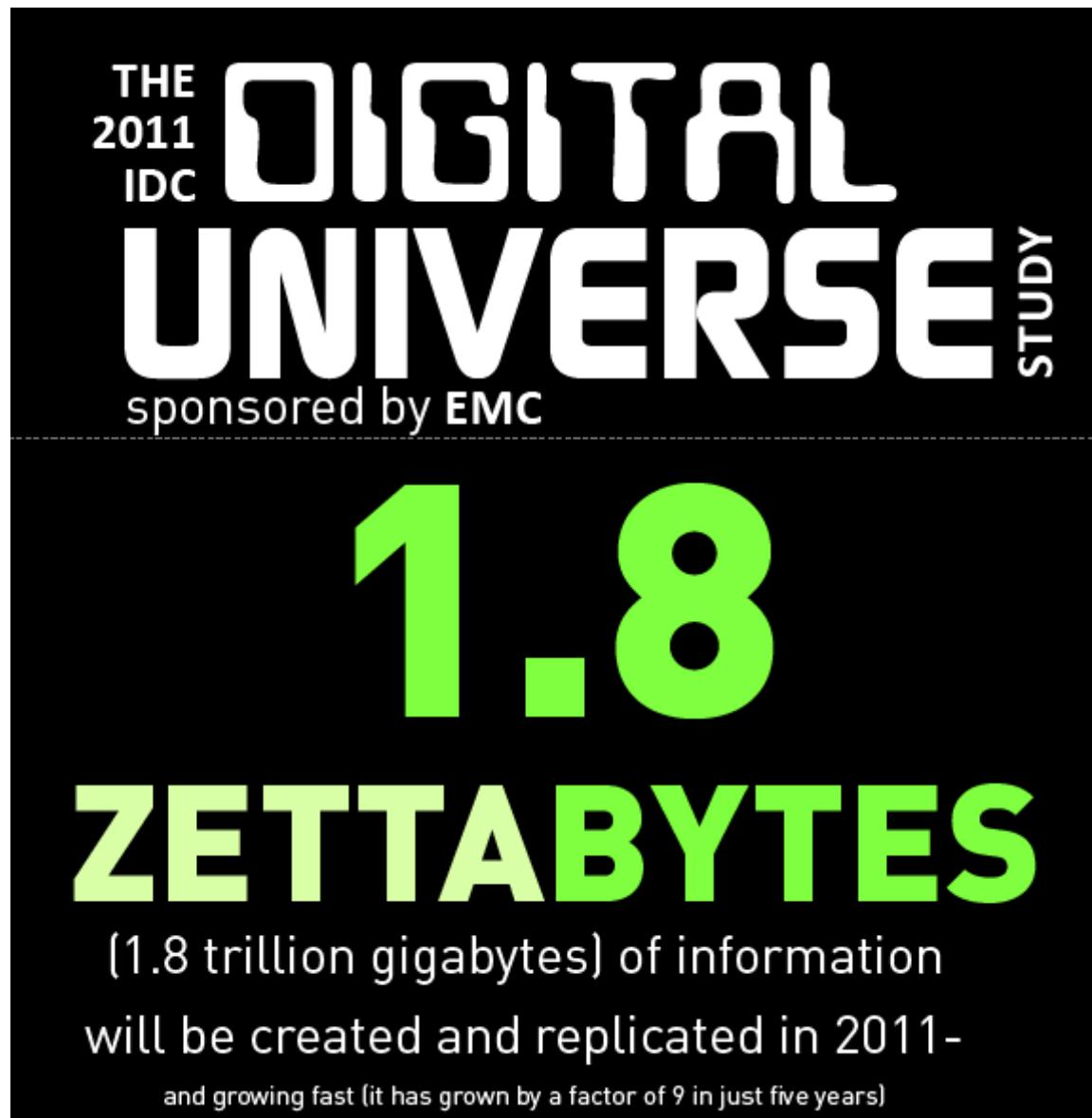
Source: IDC, 2008

2009 年，IDC 預估  
2011 年會成長十倍！  
(相較 2006 年)

Year	Volume (EB)
2006	161
2007	281
2010	988 (Forecast)
2011	1773 (Forecast)

出處：[The Diverse and Exploding Digital Universe, An Updated Forecast of Worldwide Information Growth Through 2011 March 2008, An IDC White Paper - sponsored by EMC](#)  
<http://www.emc.com/collateral/analyst-reports/diverse-exploding-digital-universe.pdf>

Data expanded 2x each year !! 每年約略兩倍



追蹤歷年的 IDC 數據：

2006	161	EB	
2007	281	EB	
2008	487	EB	
2009	800	EB	(0.8 ZB)
2010	988	EB	(預測)
2010	1200	EB	(1.2 ZB)
2011	1773	EB	(預測)
2011	1800	EB	(1.8 ZB)

景氣差而成長趨緩？  
或受新技術抑制？

出處 : Extracting Value from Chaos,

June 2011, An IDC White Paper - sponsored by EMC

<http://www.emc.com/collateral/about/news/idc-emc-digital-universe-2011-infographic.pdf>

# What is Big Data?! 何謂『海量資料』？

海量資料泛指資料大小已無法用一般軟體擷取、管理與處理；  
單一資料集大小介於數十 TB 至數 PB 的資料。

'Big Data' = few dozen TeraBytes to PetaBytes in single data set.



## Definition

[edit]

Big data is a term applied to data sets whose size is beyond the ability of commonly used software tools to capture, manage, and process the data within a tolerable elapsed time. Big data sizes are a constantly moving target currently ranging from a few dozen terabytes to many petabytes of data in a single data set.

In a 2001 research report<sup>[14]</sup> and related conference presentations, then META Group (now Gartner) analyst, Doug Laney, defined data growth challenges (and opportunities) as being three-dimensional, i.e. increasing volume (amount of data), velocity (speed of data in/out), and variety (range of data types, sources). Gartner continues to use this model for describing big data.<sup>[15]</sup>

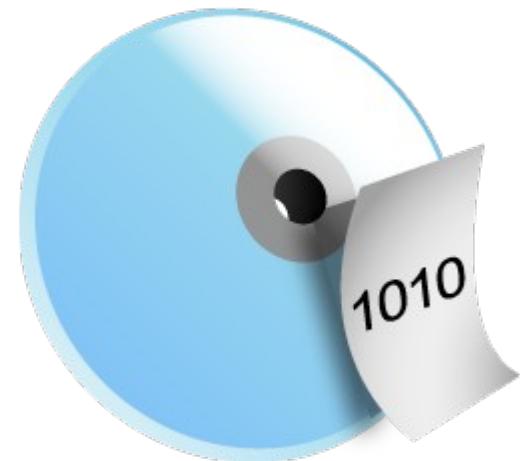
出處：[http://en.wikipedia.org/wiki/Big\\_data](http://en.wikipedia.org/wiki/Big_data)



多個檔案，容量 100TB



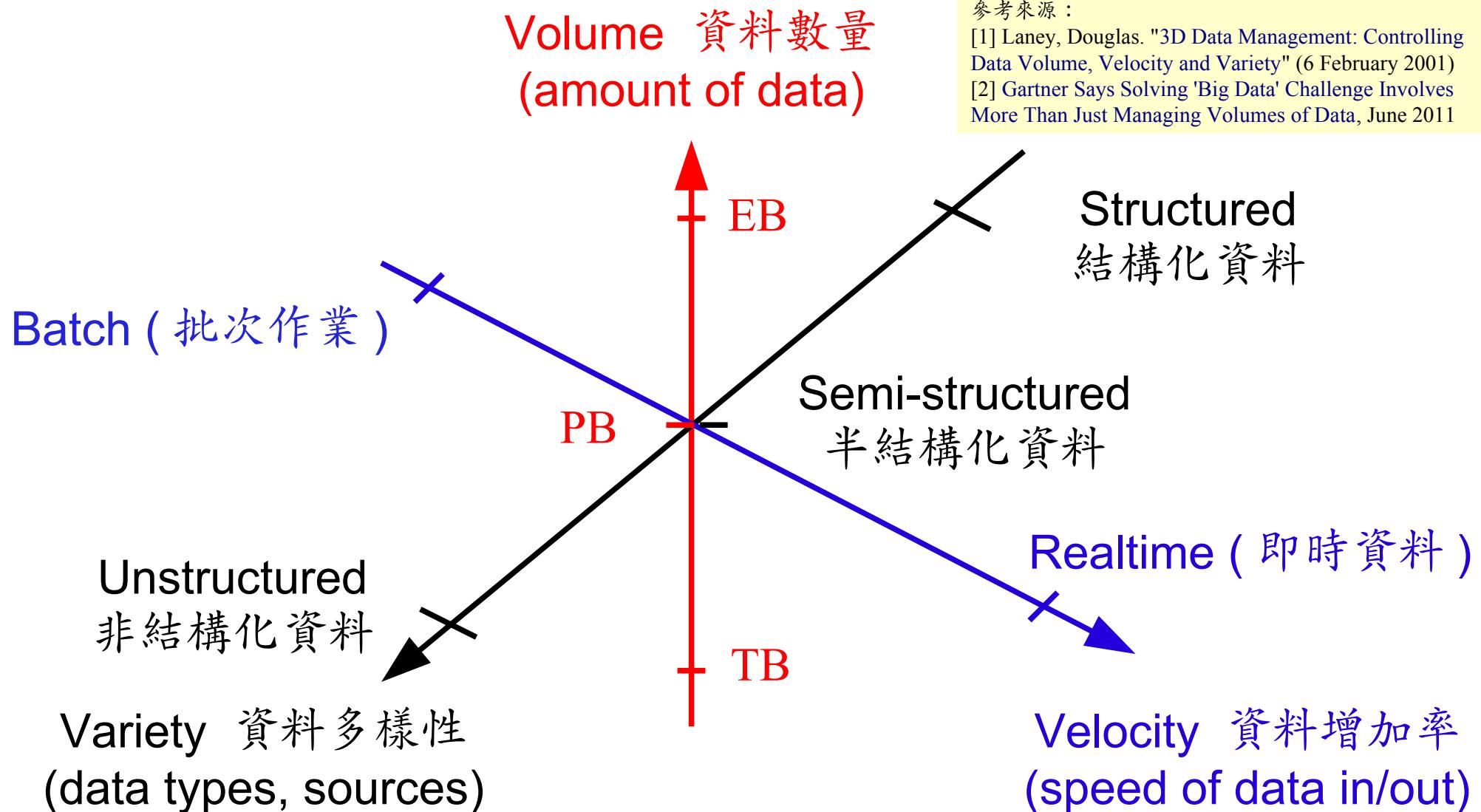
一個資料庫，容量 100TB



一個檔案，容量 100TB

# Gartner Big Data Model ? 海量資料的模型 ?

海量資料的挑戰在於如何管理「數量」、「增加率」與「多樣性」



# Six Dimensions of Big Data? 六個維度？



Source: Big Data, not Big Problems, <http://www.talend.com/products-big-data/>

# 12D of Information Management? 12 個維度？



品質管控

- Qualification and Assurance

權限管控

- Access Enablement and Control

數量管控

- Quantification

Big Data  
只是終極  
資訊管理  
的開端！

Source: Gartner (March 2011), 'Big Data' Is Only the Beginning of Extreme Information Management, 7 April 2011, <http://www.gartner.com/id=1622715>

# Agenda 演講大綱

What is Big Data ?

何謂海量資料

**Why** should we care? 為何需要關切

資料

Data

知識

Knowledge

智慧

Wisdom

**WHY**



花精灵-小麦

# Why we call it “SMART”!!

## 智慧打哪兒來？！

**Smart Phone**

智慧手機

**Smart Grid**

智慧電網

**Smart Home**

智慧家庭

**Smart Car**

智慧車輛

**Smart City**

智慧城市

**SMART**

哪裡長  
智慧了？

資料

Data

知識

Knowledge

智慧

Wisdom

# Can Machine understand You? 讓機器更懂你?

iPhone

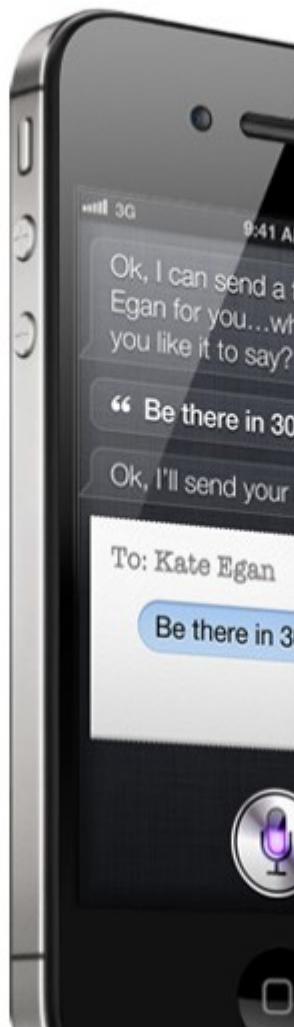
Features

Built-in Apps



Siri. Beta  
Your wish is  
its command.

Siri on iPhone 4S lets you use your voice to send messages, schedule meetings, place phone calls, and more. Ask Siri to do things just by talking the way you talk. Siri understands what you say, knows what you mean, and even talks back. Siri is so easy to use and does so much, you'll keep finding more and more ways to use it.



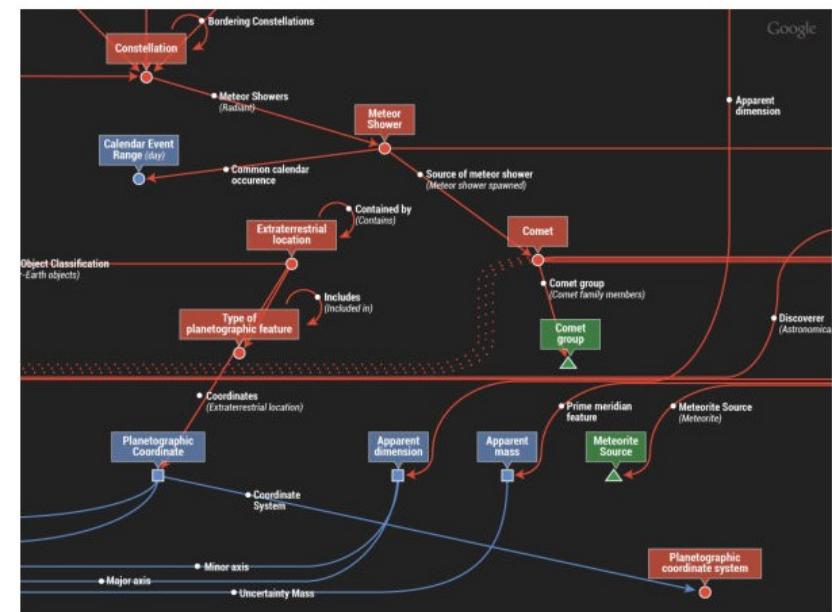
## Google將發展「人工智慧」 永久改變搜尋引擎

2012年02月15日 00:11

點評：超級阿斯拉，衝啊！（阿斯拉：好的，隼人！）

記者黃郁楨／綜合報導

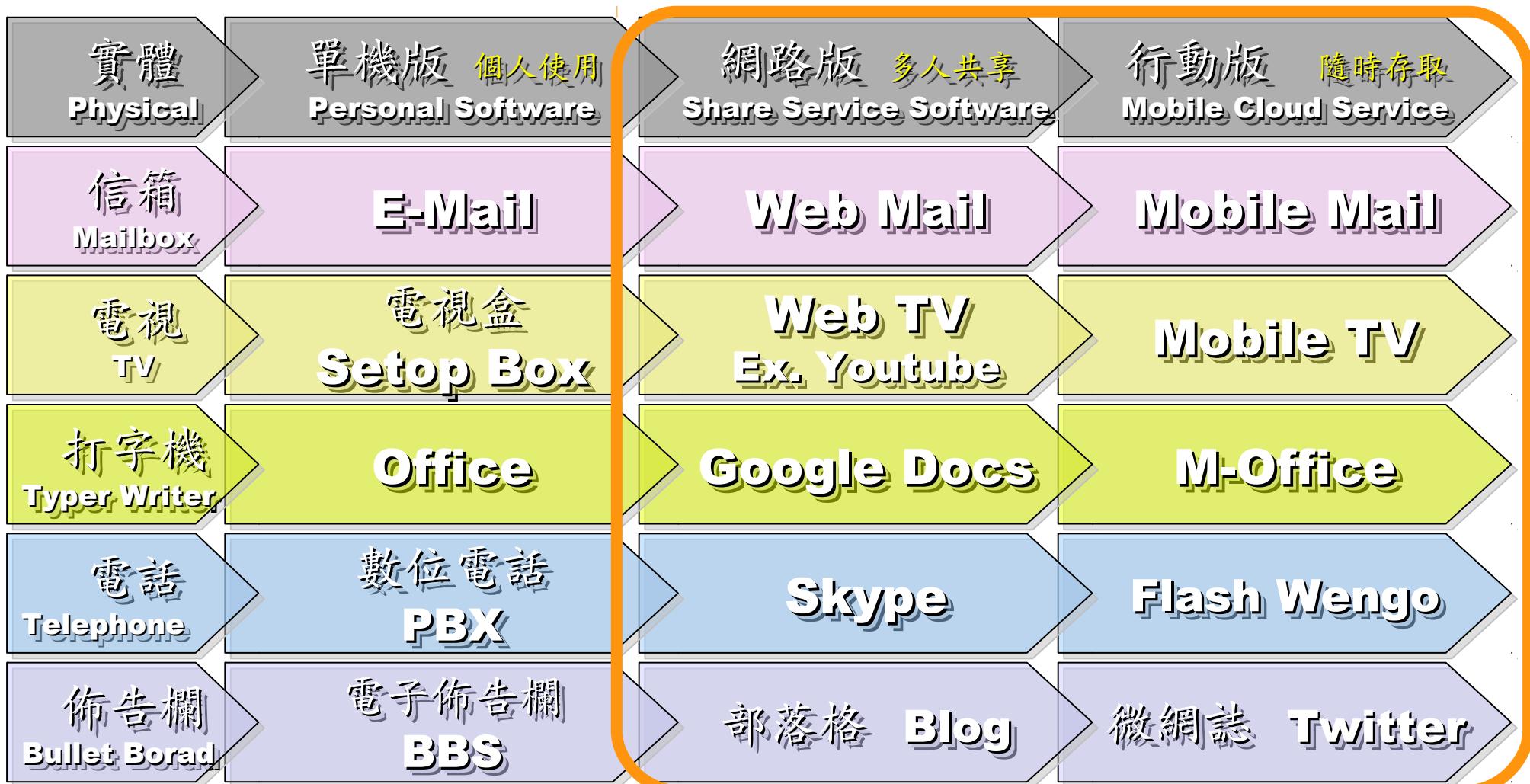
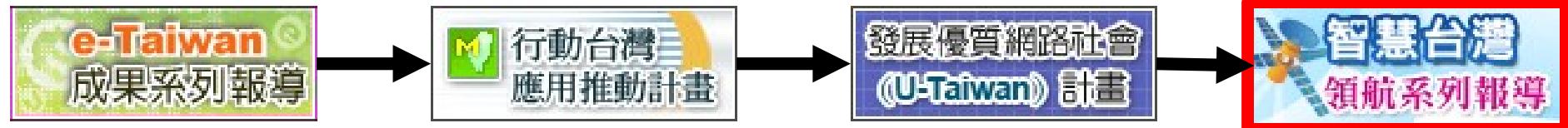
每個人都在猜，下一波網路革命是什麼？每個人都在猜，未來的世界會如何運作？Google的資深副總Amit Singhai透露了一點訊息。「Google正努力從『單字』層面進展到『意義』層面，未來搜尋引擎提供的不只是關鍵字搜尋，搜尋引擎甚至會『明白』你到底要什麼。」



▲Google未來將會朝「人工智慧」前進。（圖／取自mashable.com）

# Evolution of Software / Service

## 軟體演化勢必走向『智能化』



# The wisdom of Clouds (Crowds)

雲端序曲：雲端的智慧始終來自於群眾的智慧

2006年8月9日

Google 執行長施密特（Eric Schmidt）於SES'06會議中首次使用  
「雲端運算（Cloud Computing）」來形容無所不在的網路服務

2006年8月24日

Amazon 以 Elastic Compute Cloud 命名其虛擬運算資源服務



# Data is the source of Wisdom !!

用雲掌握資料，加以分析，形成智能給端用



雲

資料中心  
提供服務

雲端設計新思維：端的智能來自於雲的服務

**Devices share the wisdom of Cloud**

端

各類裝置  
存取服務



# Agenda 演講大綱

What is Big Data ?

何謂海量資料

Why should we care?

為何需要關切

**When to deploy it ?**

何時導入技術

基礎建設

IaaS

分析平台

PaaS

智慧服務

SaaS

**WHEN**



花精靈~小魯

# National Definition of Cloud Computing

## 美國國家標準局 NIST 給雲端運算所下的定義

### 5 Characteristics

五大基礎特徵

### 4 Deployment Models

四個佈署模型

### 3 Service Models

三個服務模式

#### 1. On-demand self-service.

隨需自助服務

#### 2. Broad network access

隨時隨地用任何網路裝置存取

#### 3. Resource pooling

多人共享資源池

#### 4. Rapid elasticity

快速重新佈署靈活度

#### 5. Measured Service

可被監控與量測的服務

# 4 Deployment Models of Cloud Computing

雲端運算的四種佈署模型

Public Cloud

公用雲端

Target Market

is **S.M.B.**

主要客戶為  
中小企業

Community Cloud

社群雲端

Academia 學術為主



**Dynamic Resource Provisioning  
between public and private cloud**

私有雲端動態根據計算需求  
調用公用雲端的資源

Hybrid  
Cloud

以大型企業  
為主要客戶  
**Enterprise** is  
key market



私有雲端  
Private Cloud

# 3 Service Models of Cloud Computing

雲端運算的三種服務模式 (市場區隔)

## IaaS

Infrastructure as a Service

架構即服務

## PaaS

Platform as a Service

平台即服務

## SaaS

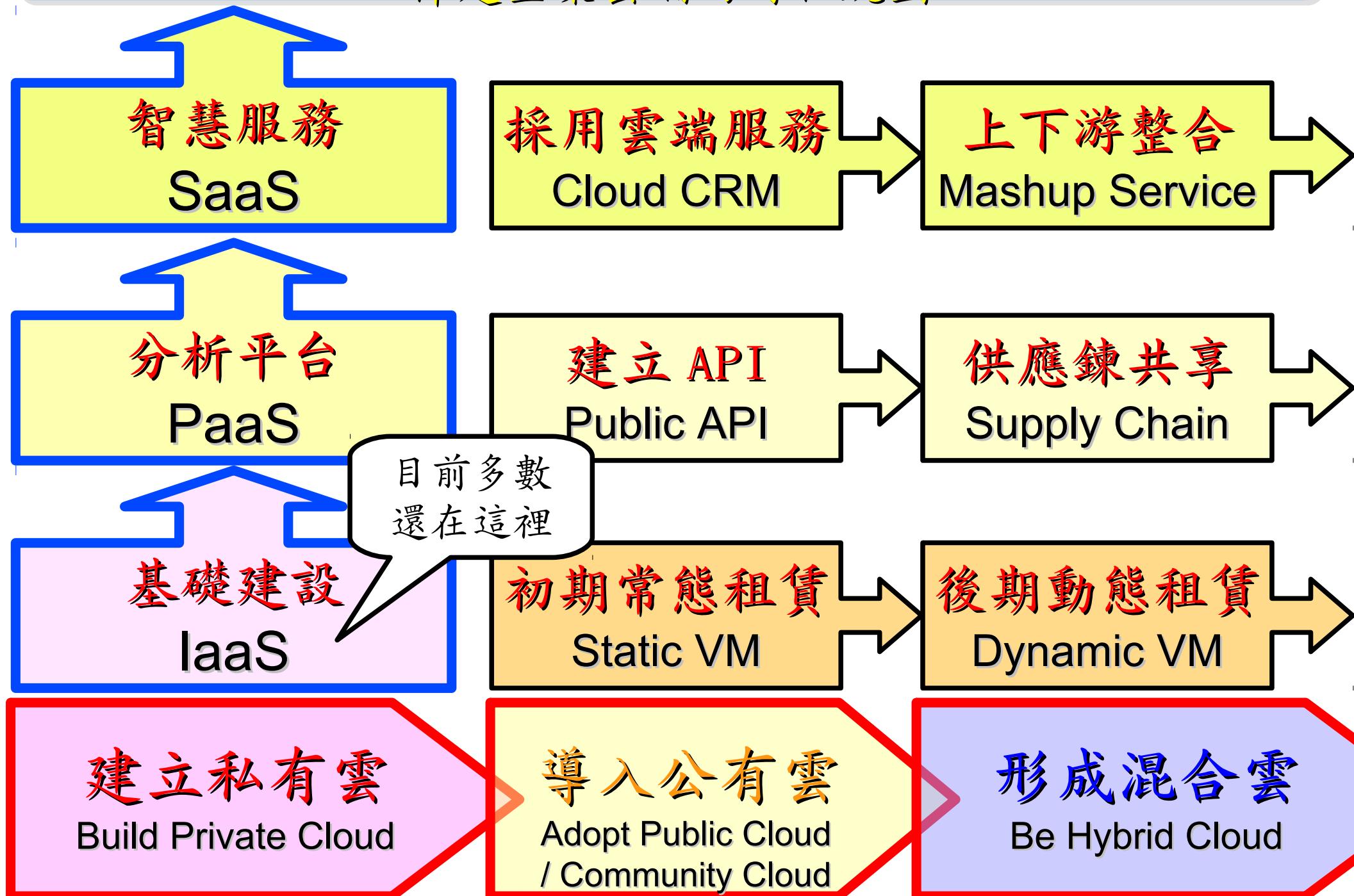
Software as a Service

軟體即服務



# Roadmap to build Your Enterprise Cloud !!

佈建企業雲端的時程規劃



# Agenda 演講大綱

**What is Big Data ?**

何謂海量資料

**Why should we care?**

為何需要關切

**When to deploy it ?**

何時導入技術

**How to handle it ?**

三大因應策略

儲存虛擬化

Dedup.

資料安全

Security

智慧服務

SaaS

**HOW**



花精靈-函兒

# Three Solutions !! 三種服務模式 vs. 三類因應對策

**SaaS**

Software as a Service

軟體即服務

**Web 2.0**

網頁服務

**PaaS**

Platform as a Service

平台即服務

**Data Analysis**

資料分析

**IaaS**

Infrastructure as a Service

架構即服務

**Virtualization**

虛擬化技術

(A) 提供 API 介面

(B) 分散式資料庫

(A) 資料整合

(B) 資料探勘

(A) 儲存虛擬化

(B) 備援與加密

# Agenda 演講大綱

**What is Big Data ?** 何謂海量資料

**Why should we care?** 為何需要關切

**When to deploy it ?** 何時導入技術

**How to handle it ?** 三大因應策略

**Who is key player ?** 誰是成功關鍵

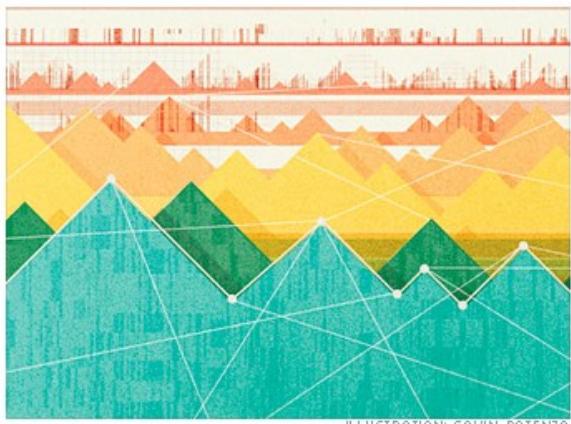


# Data Scientist !! 資料科學家 !!

## Data scientist: The hot new gig in tech

By Michal Lev-Ram, writer September 6, 2011: 5:00 AM ET

Companies that want to make sense of all their bits and bytes are hiring so-called data scientists - if they can find any.



FORTUNE -- The unemployment rate in the U.S. continues to be abysmal ([9.1% in July](#)), but the tech world has spawned a new kind of highly skilled, nerdy-cool job that companies are scrambling to fill: data scientist.

會「統計」的人照過來！  
財星雜誌 (FORTUNE) 等均報導今年最熱門的職缺是「資料科學家」！

### What is data science?

Data science can be broken down into four essential parts.

#### Mining data



Collecting and formatting  
the information

#### Statistics



Information analysis

#### Interpret



Representation or visualization in  
the form of presentations,  
infographics, graphs or charts

#### Leverage

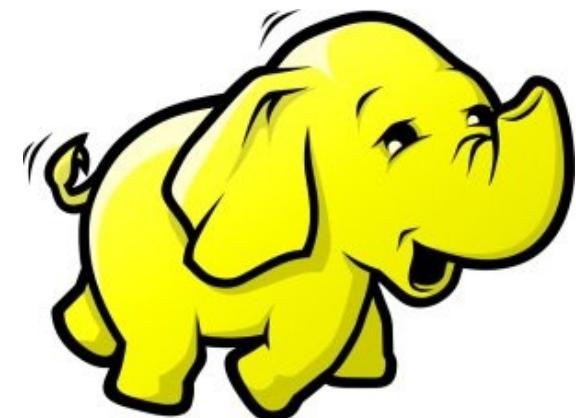


Implications of the data,  
application of the data, interaction  
using the data and predictions  
formed from studying it



## 處理海量資料的資訊架構與關鍵技術 *Technologies to build IT Stack for Big Data*

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Hot Jobs in Big Data

## 從海量資料的熱門工作談起

Data Mining

資料探勘

Data Visualization

資料視覺化

Data Analysis

資料分析

Data Manipulation

資料操控

Data Discovery

資料鑑識

How to Get a Hot Job in Big Data, Dan Tynan, InfoWorld, March 19, 2012  
出處：<http://www.cio.com/article/print/702388>

# Applications of Data Mining

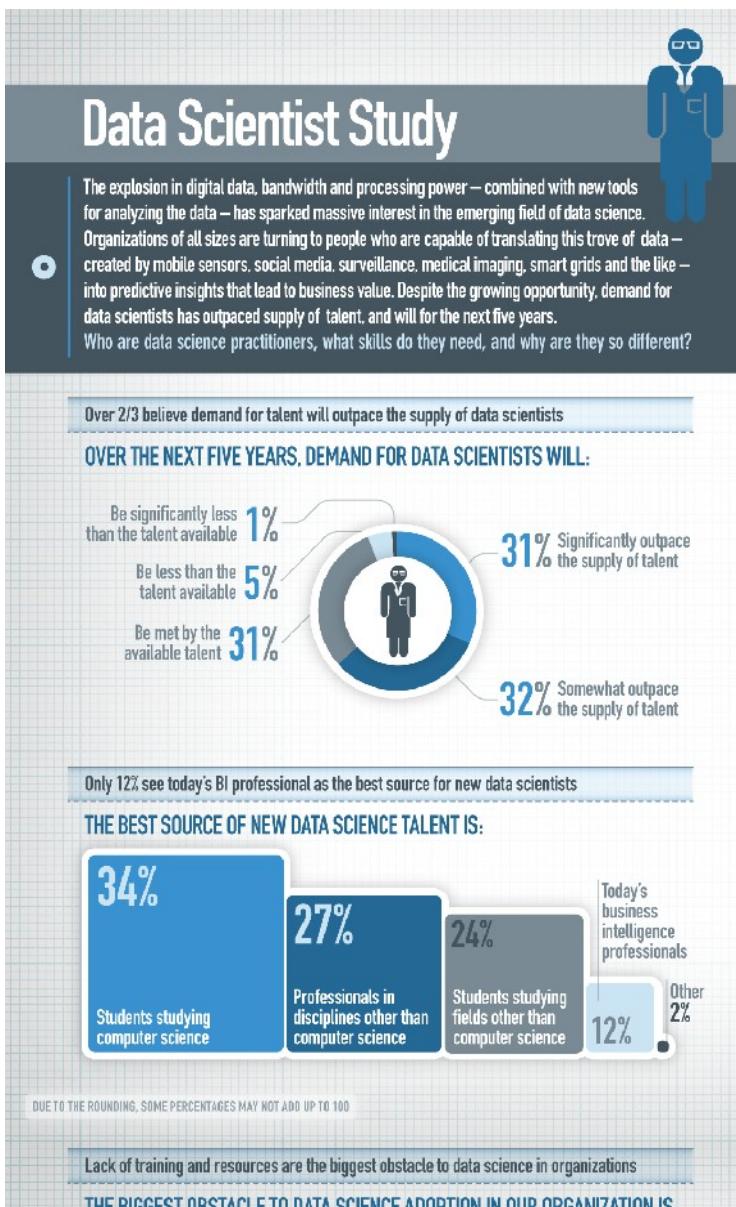
## 資料探勘的應用 - 搜尋引擎

The collage consists of five screenshots:

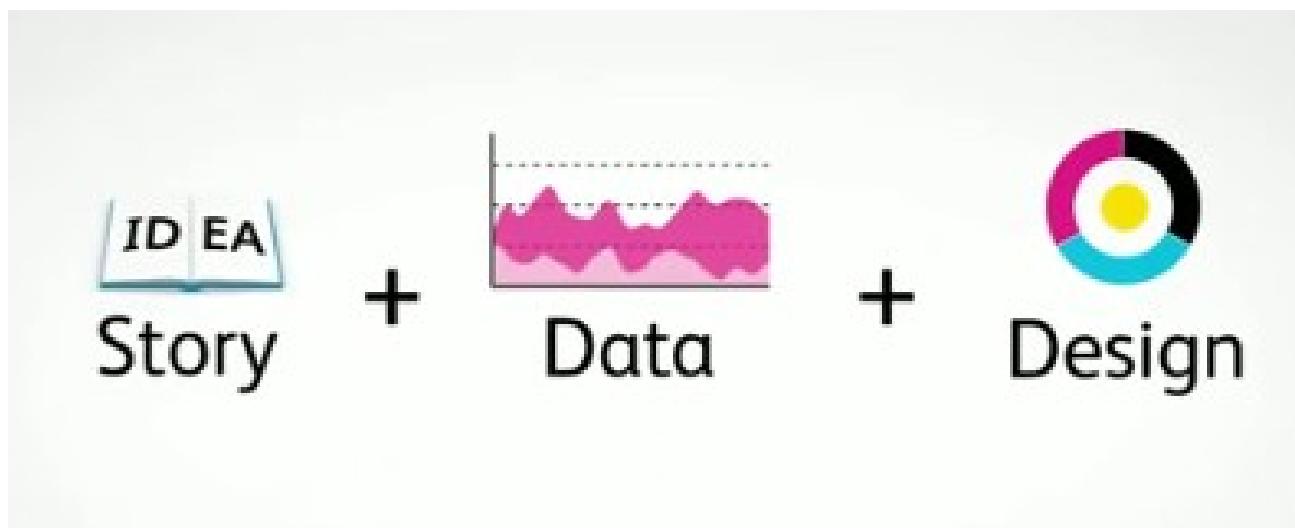
- 档案搜尋 (File Search):** A screenshot of a file search interface with a sidebar for 'Search Assistant' containing options like '圖片、音樂, 或視訊(P)', '文件 (文字處理、試算表, 等等)(Q)', and '說明和支持中心裡的資訊(I)'.
- 信件搜尋 (Email Search):** A screenshot of a Gmail search interface showing search filters for 'From', 'To', 'Subject', 'Has the words', 'Doesn't have', and 'Has attachment'.
- 即時通訊搜尋 (Instant Messaging Search):** A screenshot of a messaging interface showing a list of messages from 'jarwin.nchc.org.tw' on December 2, 2011, at 10:53:46 AM.
- 資料庫搜尋 (Database Search):** A screenshot of the IEEE Xplore Digital Library search interface, showing a search bar for 'Search 3,076,887 documents' and browse categories like 'Journals & Magazines' and 'Conference Proceedings'.
- 網頁搜尋 (Web Search):** A screenshot of the Yahoo! Taiwan homepage featuring a search bar and various search results.

Overlaid on the center of the collage are large red Chinese characters: 信件搜尋, 即時通訊搜尋, 資料庫搜尋, and 網頁搜尋.

# Applications of Data Visualization 資料視覺化的應用 - Infographics



參考來源：未來「夯」職業：資料科學家  
淺談超吸睛的資訊圖表



<http://www.bnnext.com.tw/print/article/id/21740>  
<http://www.inside.com.tw/2011/04/13/infographics>

# Applications of Data Analysis 資料分析的應用 - 商業智慧 (BI)



# Applications of Data Discovery

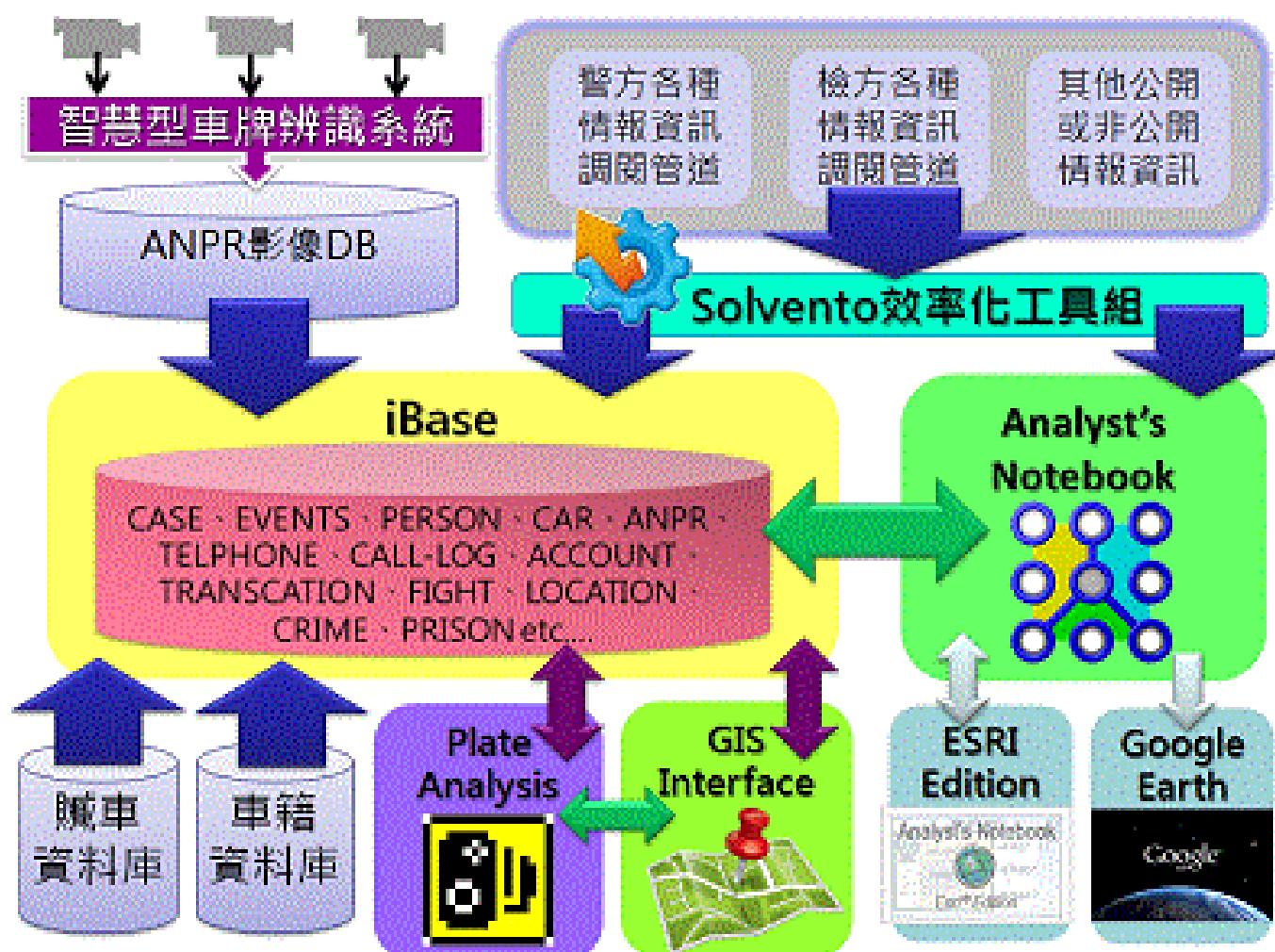
## 數位鑑識 - 資訊與法律的結合

### 電腦鑑識&會計鑑識



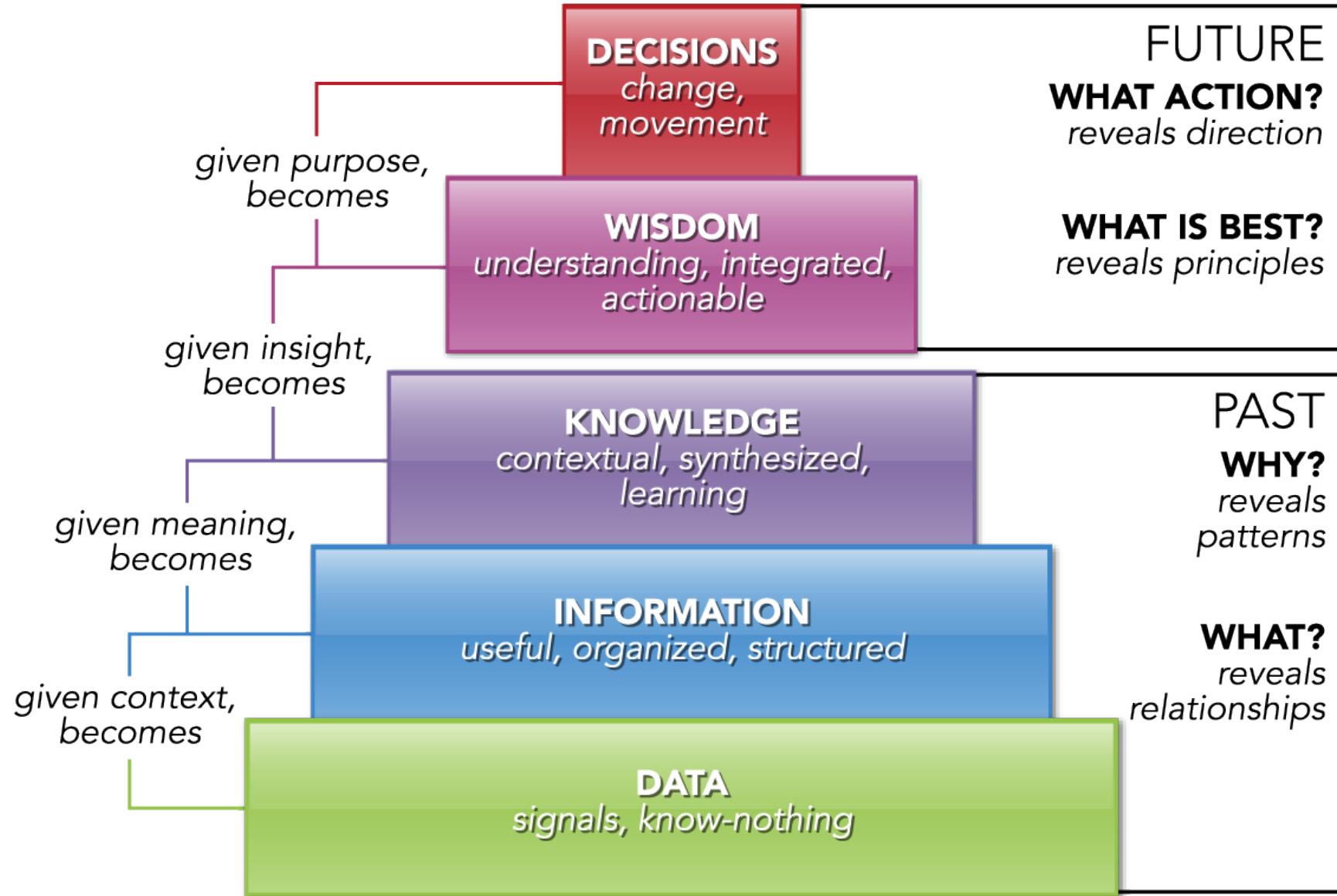
<http://blog.udn.com/kf0630/6018593>

[http://www.solventosoft.com/upload/ANPR\\_02s.gif](http://www.solventosoft.com/upload/ANPR_02s.gif)



# Data, Information, Knowledge, Wisdom

## 知識管理模型：資料、資訊、知識與智慧



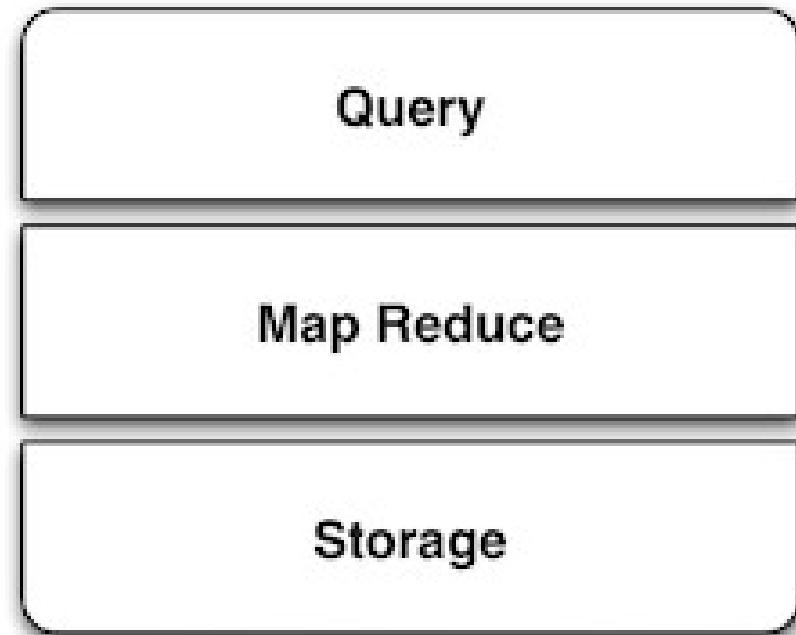
# The SMAQ stack for big data

## 海量資料處理的資訊架構

做網頁相關的人可能聽過 LAMP



未來處理海量資料的人必需知道  
SMAQ ( Storage, MapReduce and Query )

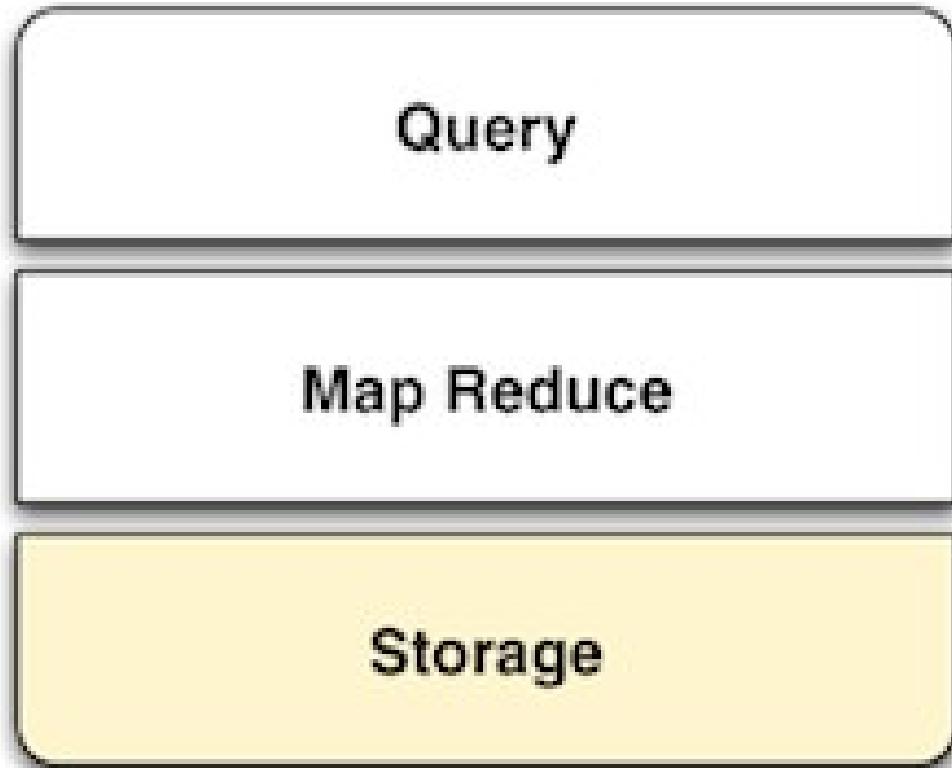


參考來源：The SMAQ stack for big data , Edd Dumbill , 22 September 2010 ,  
<http://radar.oreilly.com/2010/09/the-smaq-stack-for-big-data.html>

圖片來源：<http://smashingweb.ge6.org/wp-content/uploads/2011/10/apache-php-mysql-ubuntu.png> 37

# The SMAQ stack for big data

## 海量資料處理的資訊架構

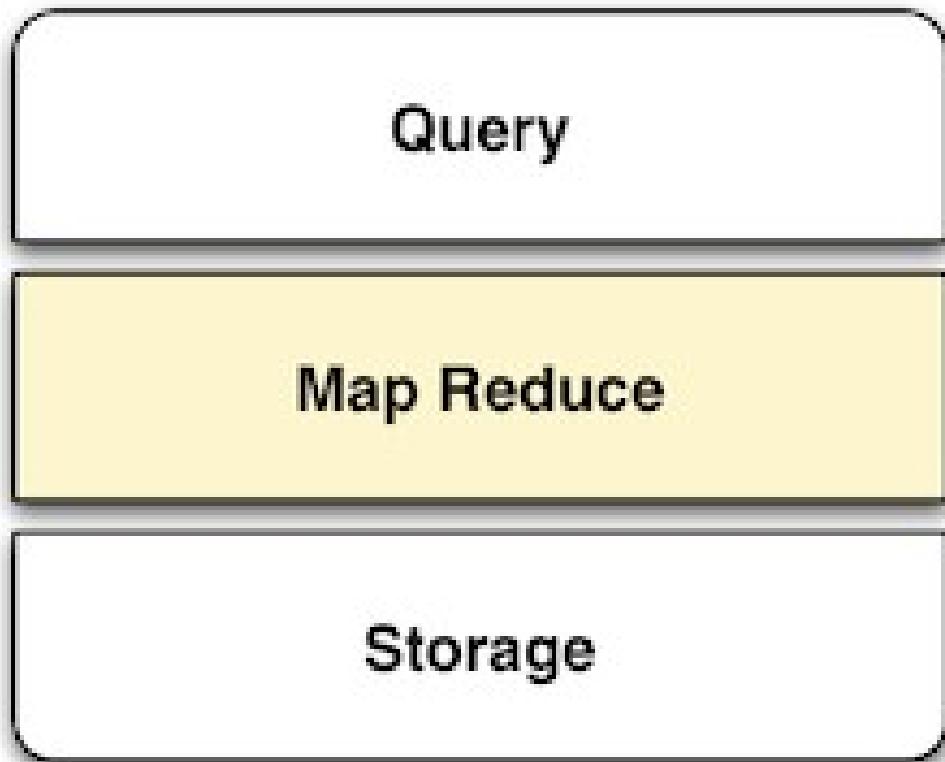


用來儲存分散、沒有關聯  
的非結構化資料



# The SMAQ stack for big data

## 海量資料處理的資訊架構



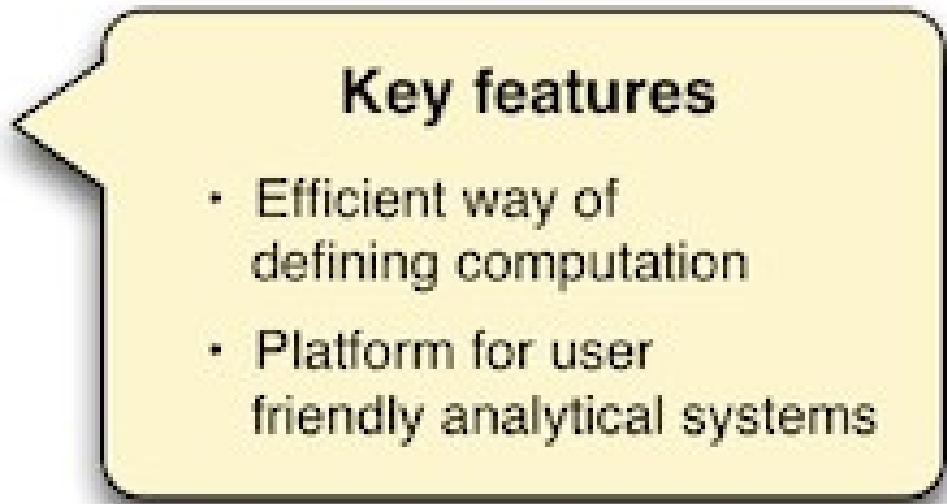
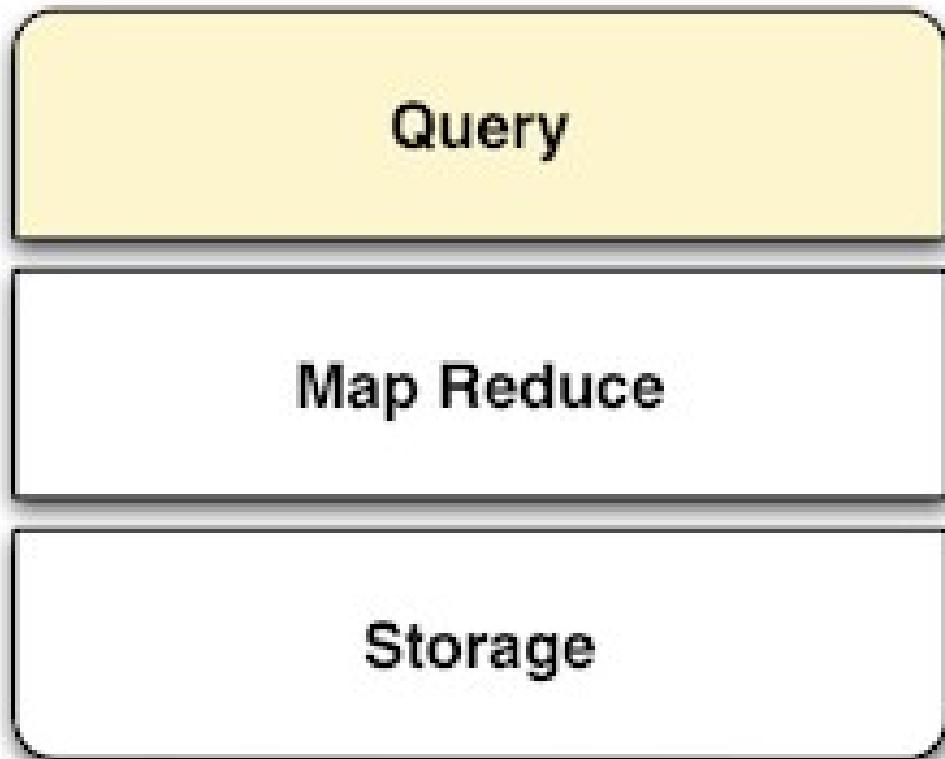
運用批次處理的方式，將運算工作平均分散到許多的伺服器做運算。

### Key features

- Distributes computation over many servers
- Batch processing model

# The SMAQ stack for big data

## 海量資料處理的資訊架構



將算完的結構化資料儲存到可供查詢的資料庫系統

# Three Core Technologies of Google ....

## Google 的三大關鍵技術 ....

- Google 在一些會議分享他們的三大關鍵技術
- Google shared their design of web-search engine
  - SOSP 2003 :
  - “The Google File System”
  - <http://labs.google.com/papers/gfs.html>
- OSDI 2004 :
  - “MapReduce : Simplified Data Processing on Large Cluster”
  - <http://labs.google.com/papers/mapreduce.html>
- OSDI 2006 :
  - “Bigtable: A Distributed Storage System for Structured Data”
  - <http://labs.google.com/papers/bigtable-osdi06.pdf>



# Open Source Mapping of Google Core Technologies

## Google 三大關鍵技術對應的自由軟體

### BigTable

A huge key-value datastore

### HBase, Hypertable

Cassandra, ....

### MapReduce

To parallel process data

### Hadoop MapReduce API

Sphere MapReduce API, ...

### Google File System

To store petabytes of data

### Hadoop Distributed File System (HDFS)

Sector Distributed File System

更多不同語言的 MapReduce API 實作：

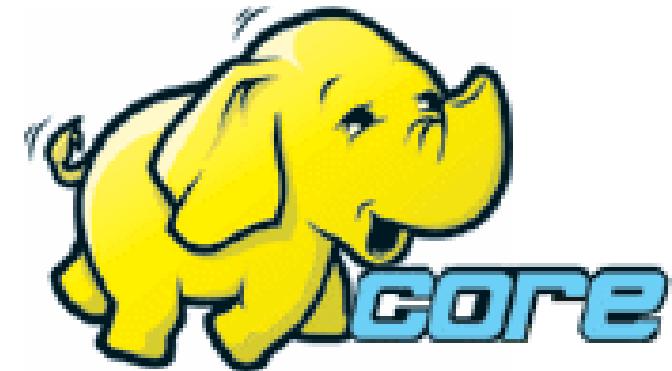
<http://trac.nchc.org.tw/grid/intertrac/wiki%3Ajazz/09-04-14%23MapReduce>

其他值得觀察的分散式檔案系統：

- IBM GPFS - <http://www-03.ibm.com/systems/software/gpfs/>
- Lustre - <http://www.lustre.org/>
- Ceph - <http://ceph.newdream.net/>

# Hadoop

- <http://hadoop.apache.org>
- Hadoop 是 Apache Top Level 開發專案
- **Hadoop is Apache Top Level Project**
- 目前主要由 Yahoo! 資助、開發與運用
- **Major sponsor is Yahoo!**
- 創始者是 Doug Cutting ，參考 Google Filesystem
- **Developed by Doug Cutting, Reference from Google Filesystem**
- 以 Java 開發，提供 HDFS 與 MapReduce API 。
- **Written by Java, it provides HDFS and MapReduce API**
- 2006 年使用在 Yahoo 內部服務中
- **Used in Yahoo since year 2006**
- 已佈署於上千個節點 。
- **It had been deploy to 4000+ nodes in Yahoo**
- 處理 Petabyte 等級資料量 。
- **Design to process dataset in Petabyte**



Facebook、Last.fm  
、Joost are also  
powered by Hadoop

# Sector / Sphere

- <http://sector.sourceforge.net/>
- 由美國資料探勘中心研發的自由軟體專案。
- Developed by National Center for Data Mining, USA
- 採用 C/C++ 語言撰寫，因此效能較 Hadoop 更好。
- Written by C/C++, so performance is better than Hadoop
- 提供「類似」Google File System 與 MapReduce 的機制
- Provide file system similar to Google File System and MapReduce API
- 基於UDT高效率網路協定來加速資料傳輸效率
- Based on UDT which enhance the network performance
- Open Cloud Testbed有提供測試環境，並開發MalStone效能評比軟體
- Open Cloud Consortium provide Open Cloud Testbed and develop MalStone toolkit for benchmark



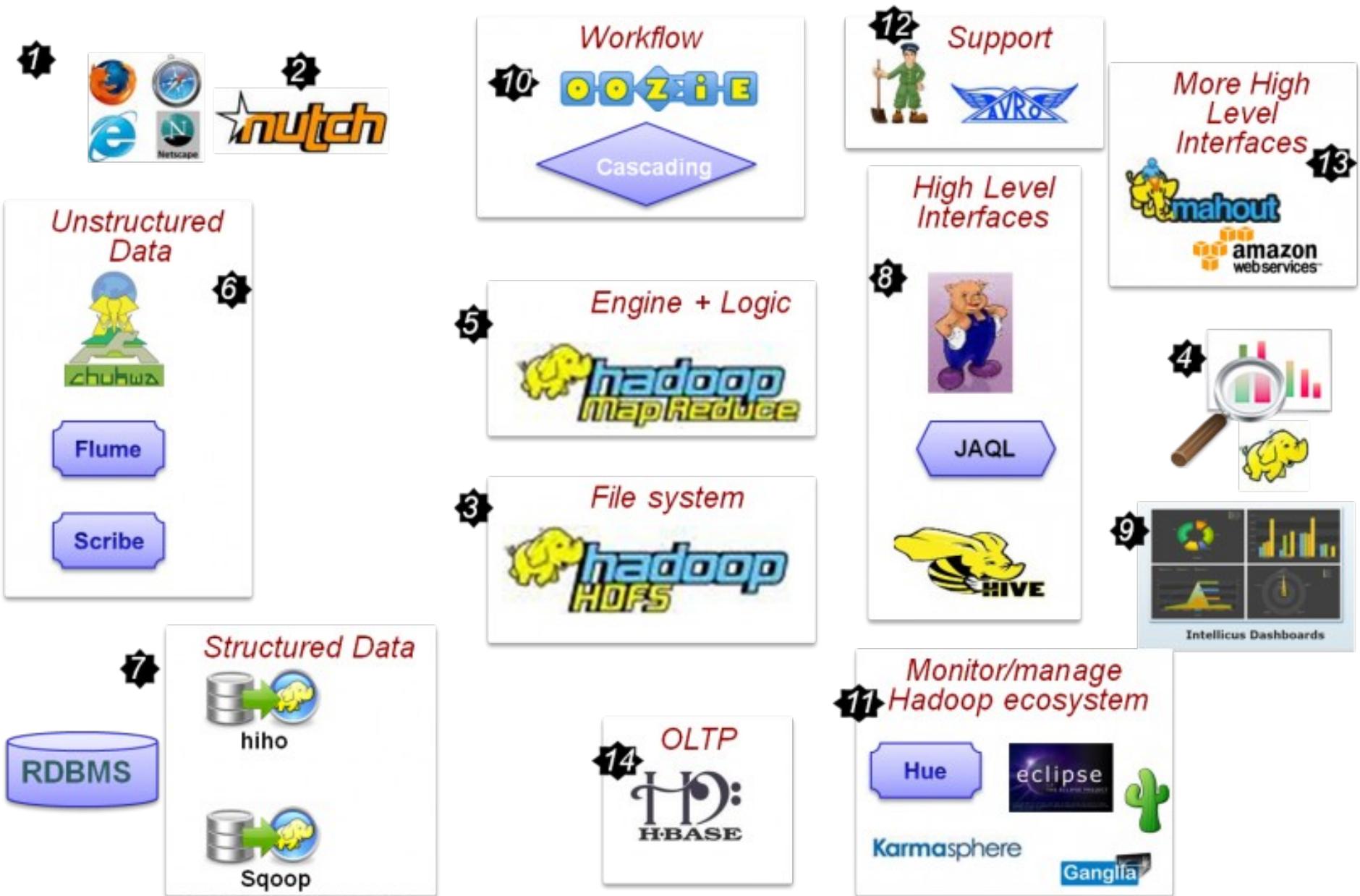
National Center for Data Mining  
University of Illinois at Chicago



Open Data Group  
<http://www.opendatagroup.com/>

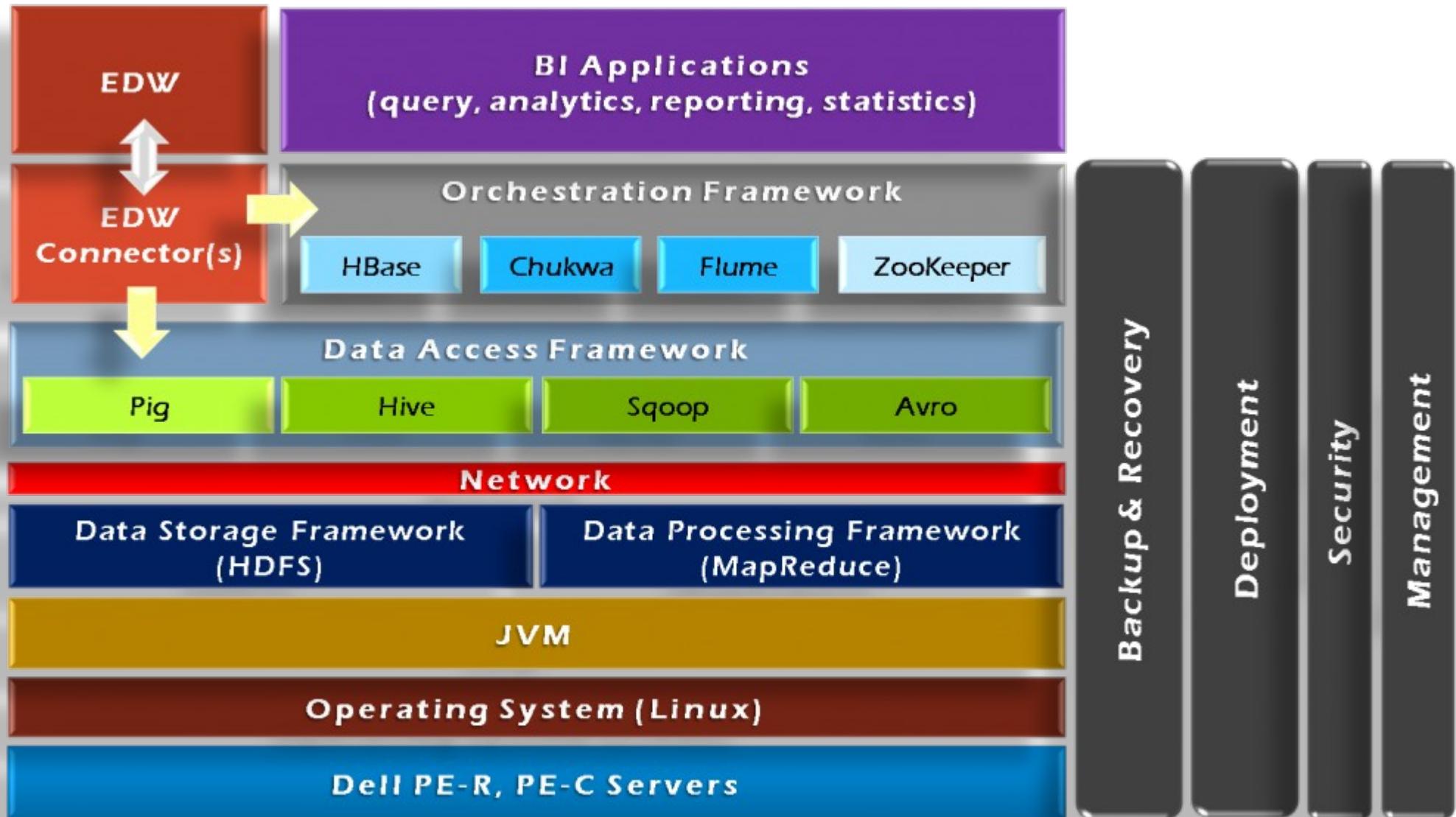
# Why we choice Hadoop? Good Ecosystem!

豐富的生態系建構出處理海量資料的工具庫



# BI and EDW build on Hadoop Ecosystem

## 運用 Hadoop 生態系搭建資料倉儲與商業智慧分析



# Build your own search engine, too

您也能用 **Hadoop** 搭建自己的搜尋引擎

Web UI ( Crawlzilla Website + Search Engine)

JSP + Servlet + JavaBean

Nutch

Lucene

Crawlzilla System Management

Tomcat

Hadoop

PC1

PC2

PC3

# Microsoft love Hadoop, too 微軟幫 Azure 還有 SQL Server 都接上 Hadoop

The screenshot shows the Microsoft SQL Server website's Business Intelligence section. At the top, there are navigation links for About SQL Server, Solutions & Technologies, Editions, Get SQL Server, Learning Center, and Partners. Below the navigation is a red banner with the text "Big Data Analytics". To the left, there is a video player for a video titled "Strata Big Data Conference 2012 and Power View Contest". The video player interface includes a play button, volume controls, a progress bar showing 0:00 / 2:35, and a YouTube logo. To the right of the video player is a section titled "Big Data Solution" which discusses unlocking business insights from structured and unstructured data using Microsoft's Big Data solution, which integrates Apache Hadoop into various Microsoft BI tools like SQL Server Analysis Services and Reporting Services.

## Big Data Solution

Unlock business insights from all your structured and unstructured data, including large volumes of data not previously activated, with Microsoft's Big Data solution. Microsoft's end-to-end roadmap for Big Data embraces Apache Hadoop™ by distributing enterprise class Hadoop based solutions on both Windows Server and Windows Azure. Our solution is also integrated into the Microsoft BI tools such as SQL Server Analysis Services, Reporting Services and even PowerPivot and Excel. This enables you to do BI on all your data, including those in Hadoop.

### Key Benefits

- Broader access of Hadoop to end users, IT professionals and Developers, through easy installation and configuration and simplified programming with JavaScript.
- Enterprise ready Hadoop distribution with greater security, performance, ease of management and options for Hybrid IT usage.

參考來源：Big Data Solution | Microsoft SQL Server 2008 R2

<http://www.microsoft.com/sqlserver/en/us/solutions-technologies/business-intelligence/big-data-solution.aspx>

# Oracle love Hadoop, too Oracle 也接上 Hadoop

The screenshot shows the CNET News homepage with a banner for HP Officejet Pro printers. Below the banner, a news article is displayed with the title "Cloudera teams up to connect Oracle and Hadoop".

CNET › News › Software, Interrupted

## Cloudera teams up to connect Oracle and Hadoop

Cloudera and Quest software are partnering to provide connectivity between Oracle and Hadoop.



by [Dave Rosenberg](#) | June 21, 2010 5:30 AM PDT

[Follow](#)

This week [Cloudera](#), a provider of software and services for the Apache Hadoop project, is set to announce a partnership with [Quest Software](#) to develop, support, and distribute an Oracle connector for Hadoop.



參考來源 : Cloudera teams up to connect Oracle and Hadoop

[http://news.cnet.com/8301-13846\\_3-20008242-62.html](http://news.cnet.com/8301-13846_3-20008242-62.html)

# Hinet Application of Big Data

## 中華電信已經在做的海量資料應用



### 中華電信：分析駭客行為，拓展對外新服務

撰文者：趙郁竹

發表日期：2012-03-06



#### [214期雜誌精選]

全台最大的中華電信提供行動電話、市話、寬頻固網、MOD……，各種業務服務，加起來的用戶數就有3000萬，比全台灣人口還多，光是單月帳務數量就高達100億筆資料。除了電信、寬頻服務，還有日益增加的數位服務、行動增值服務，從服務內容到客戶端，累積出的資料相當驚人。

「資料量越來越大，日常分析工作需要很多時間，但新的運算技術有效解決了這個問題，」中華電信資訊處處長陳明仕說。2010年開始，因為中華電信本身的資料運算需求，採用分散式運算架構Hadoop技術，打造出大資料運算平台，不但解決了自身的資料問題，還能對外提供資料運算應用。

以MOD為例，一天有幾千萬筆資料，如何找出使用者在什麼時段做了什麼事？廣告效益又如何？「用傳統的方法，需要400分鐘才能分析完；用Hadoop大資料平台，13分鐘就能解決，節省非常多時間，」他說。

#### 追蹤再拆解

大資料運算技術除了節省時間，還能防止駭客入侵。「駭客的攻擊行為都有模式可循，」陳明仕解釋，就像球賽一樣，了解進攻模式就能防守。用戶的資料保護是第一要務，因此透過行為模式分析，能有效保護企業資訊安全，也保障客戶的個資安全。

參考來源：中華電信：分析駭客行為，拓展對外新服務，發表日期：2012-03-06

<http://www.bnnext.com.tw/print/article/id/22333>

# Hinet Application of Big Data

## 中華電信已經在做的海量資料應用

IT ithome.com.tw

### 中華電信用Hadoop技術分析通話明細

 READ LATER

面對資料快速成長以及非結構性資料的增加，中華電信資訊處第四科科長楊秀一表示，中華電信近來利用Hadoop雲端運算技術自行開發了一個專門用來分析非結構化資料的巨量資料（Big Data）運算平臺，嘗試在資料進到資料倉儲系統之前，先進行資料的分析與處理以減少資料倉儲的資料量。

近年來行動語音市場趨於飽和，為了掌握用戶特性進行客製化行銷，一份資料要進行分析，就會被多次複製，因此即使用戶增加趨緩，但中華電信擁有的資料量仍快速暴增。

中華電信用來分析的資料模型最早於10多年前已有雛形，但當初主要用於行動語音分析。一直到2009年，他們完整導入Teradata的電信業邏輯資料模型cLDM 9.0版，整合更多電信服務的用戶資料。楊秀一表示，當初導入該模型的目的主要是為了整合行動語音、固網、數據的資料，進行以人為中心的分析模式。在導入之前，中華電信的資料模型是以設備為中心，因為不同設備的記錄資料儲存在不同的資料庫，無法進行整合性的分析。

參考來源：中華電信用 Hadoop 技術分析通話明細，發表日期：2011-06-12

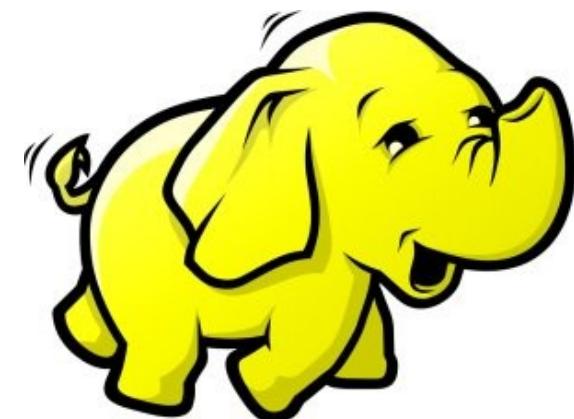
<http://www.ithome.com.tw/itadm/article.php?c=68023>



# Hadoop 簡介：源起與術語

Introduction to Hadoop : History and Terminology

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# What is Hadoop ?

用一句話解釋 **Hadoop** 是什麼 ??

*Hadoop is a **software platform**  
that lets one easily write and run  
applications that **process vast**  
**amounts of data.***

**Hadoop** 是一個讓使用者簡易撰寫並執行處理海量資料應用程式的軟體平台。

亦可以想像成一個處理海量資料的生產線，只須學會定義 **map** 跟 **reduce** 工工作站該做哪些事情。

# Features of Hadoop ...

## **Hadoop** 這套軟體的特色是 ...

- 海量 **Vast Amounts of Data**
  - 擁有儲存與處理大量資料的能力
  - Capability to STORE and PROCESS vast amounts of data.
- 經濟 **Cost Efficiency**
  - 可以用在由一般 PC 所架設的叢集環境內
  - Based on large clusters built of commodity hardware.
- 效率 **Parallel Performance**
  - 透過分散式檔案系統的幫助，以致得到快速的回應
  - With the help of HDFS, Hadoop have better performance.
- 可靠 **Robustness**
  - 當某節點發生錯誤，能即時自動取得備份資料及佈署運算資源
  - Robustness to add and remove computing and storage resource without shutdown entire system.

# Founder of Hadoop – Doug Cutting

**Hadoop** 這套軟體的創辦人 **Doug Cutting**

Doug Cutting Talks About The Founding Of Hadoop

clouderahadoop

9 部影片

編輯訂閱項目



Doug Cutting Talks About The Founding Of Hadoop

<http://www.youtube.com/watch?v=qxC4urJOchs>

# History of Hadoop ... 2002~2004

**Hadoop** 這套軟體的歷史源起 ... 2002~2004



- Lucene

- <http://lucene.apache.org/>
  - 用Java 設計的高效能文件索引引擎 API
  - a high-performance, full-featured **text search engine library** written entirely in **Java**.
  - 索引文件中的每一字，讓搜尋的效率比傳統逐字比較還要高的多
  - Lucene create an **inverse index** of every word in different documents. It enhance performance of text searching.

# History of Hadoop ... 2002~2004

**Hadoop** 這套軟體的歷史源起 ... 2002~2004

- Nutch



- <http://nutch.apache.org/>
- Nutch 是基於開放原始碼所開發的網站搜尋引擎
- Nutch is open source web-search software.
- 利用 Lucene 函式庫開發
- It builds on Lucene and Solr, adding web-specifics, such as a crawler, a link-graph database, parsers for HTML and other document formats, etc.



# Three Gifts from Google ....

## 來自 **Google** 的三個禮物 ....

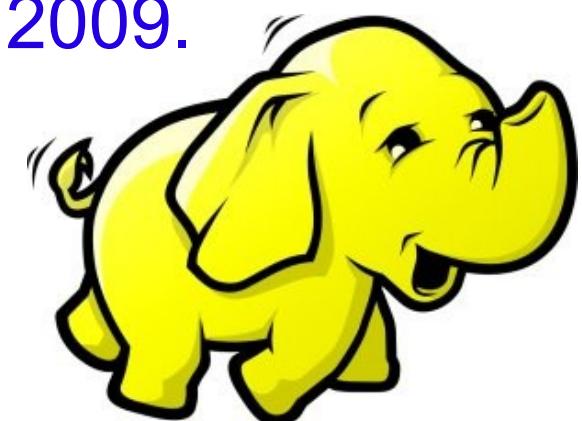
- Nutch 後來遇到儲存大量網站資料的瓶頸
- Nutch encounter storage issue
- Google 在一些會議分享他們的三大關鍵技術
- Google shared their design of web-search engine
  - SOSP 2003 : “The Google File System”
  - <http://labs.google.com/papers/gfs.html>
  - OSDI 2004 : “MapReduce : Simplified Data Processing on Large Cluster”
  - <http://labs.google.com/papers/mapreduce.html>
  - OSDI 2006 : “Bigtable: A Distributed Storage System for Structured Data”
  - <http://labs.google.com/papers/bigtable-osdi06.pdf>



# History of Hadoop ... 2004 ~ Now

## Hadoop 這套軟體的歷史源起 ... 2004 ~ Now

- Dong Cutting reference from Google's publication
- Added DFS & MapReduce implement to Nutch
- According to **user feedback** on the mail list of Nutch ....
- Hadoop became separated project **since Nutch 0.8**
- Nutch DFS → Hadoop Distributed File System (HDFS)
- **Yahoo** hire Dong Cutting to build a team of web search engine at **year 2006**.
  - Only **14 team members** (engineers, clusters, users, etc.)
- Doung Cutting joined Cloudera at year 2009.



# Who Use Hadoop ??

有哪些公司在用 **Hadoop** 這套軟體 ??

- Yahoo is the key contributor currently.
- IBM and Google teach Hadoop in universities ...
- [http://www.google.com/intl/en/press/pressrel/20071008\\_ibm\\_univ.html](http://www.google.com/intl/en/press/pressrel/20071008_ibm_univ.html)
- The New York Times used 100 Amazon EC2 instances and a Hadoop application to process 4TB of raw image TIFF data (stored in S3) into 11 million finished PDFs in the space of 24 hours at a computation cost of about \$240 (not including bandwidth)
  - from <http://en.wikipedia.org/wiki/Hadoop>
- <http://wiki.apache.org/hadoop/AmazonEC2>
- <http://wiki.apache.org/hadoop/PoweredBy>

- A9.com
- ADSDAQ by Contextweb
- EHarmony
- Facebook
- Fox Interactive Media

- IBM
- ImageShack
- ISI
- Joost
- Last.fm

- Powerset
- The New York Times
- Rackspace
- Veoh
- Metaweb

# Hadoop in production run ....

## 商業運轉中的 **Hadoop** 應用 ....

- February 19, 2008
- Yahoo! Launches World's Largest Hadoop Production Application
- <http://developer.yahoo.net/blogs/hadoop/2008/02/yahoo-worlds-largest-production-hadoop.html>

Number of links between pages in the index	roughly 1 trillion links
Size of output	over 300 TB, compressed!
Number of cores used to run single Map-Reduce job	over 10,000
Raw disk used in the production cluster	over 5 Petabytes

# Hadoop in production run ....

## 商業運轉中的 **Hadoop** 應用 ....

- September 30, 2008
- Scaling Hadoop to 4000 nodes at Yahoo!
- [http://developer.yahoo.net/blogs/hadoop/2008/09/scaling\\_hadoop\\_to\\_4000\\_nodes\\_a.html](http://developer.yahoo.net/blogs/hadoop/2008/09/scaling_hadoop_to_4000_nodes_a.html)

<b>Total Nodes</b>	<b>4000</b>
<b>Total cores</b>	<b>30000</b>
<b>Data</b>	<b>16PB</b>

	<b>500-node cluster</b>		<b>4000-node cluster</b>	
	<b>write</b>	<b>read</b>	<b>write</b>	<b>read</b>
<b>number of files</b>	<b>990</b>	<b>990</b>	<b>14,000</b>	<b>14,000</b>
<b>file size (MB)</b>	<b>320</b>	<b>320</b>	<b>360</b>	<b>360</b>
<b>total MB processes</b>	<b>316,800</b>	<b>316,800</b>	<b>5,040,000</b>	<b>5,040,000</b>
<b>tasks per node</b>	<b>2</b>	<b>2</b>	<b>4</b>	<b>4</b>
<b>avg. throughput (MB/s)</b>	<b>5.8</b>	<b>18</b>	<b>40</b>	<b>66</b>

# Comparison between Google and Hadoop

## *Google* 與 *Hadoop* 的比較表

Develop Group	Google	Apache
Sponsor	Google	Yahoo, Amazon
Algorithm Method	MapReduce	MapReduce
Resource	open document	open source
File System (MapReduce)	GFS	HDFS
Storage System (for structure data)	big-table	HBase
Search Engine	Google	Nutch
OS	Linux	Linux / GPL

# Why should we learn Hadoop ?

## 爲何需要學習 **Hadoop ??**

[Search Jobs](#) [Browse Jobs](#) [Local Jobs](#) [Salaries](#) [Employment Trends](#)



Employment Trends

Xen, Hyper-V, Hadoop

Tip: You can compare trends by separating them with commas.

Xen, Hyper-v, Hadoop Trends



### Xen, Hyper-v, Hadoop Job Trends

This graph displays the percentage of jobs with your search terms anywhere in the job listing. Since November 2008, the following has occurred:

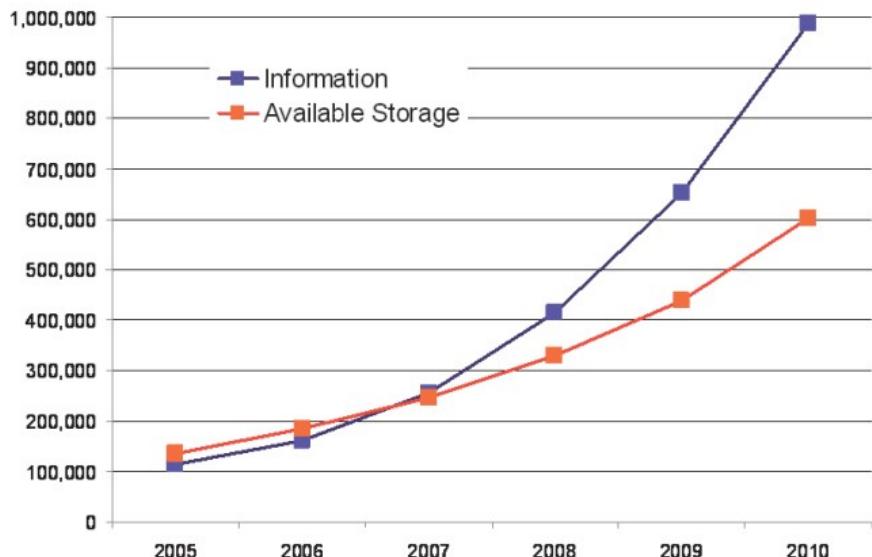
- Xen jobs increased 141%
- Hyper-v jobs increased 551%
- Hadoop jobs did not change or there is no data available

1. Data Explore  
資訊大爆炸

2. Data Mining Tool  
方便作資料探勘的工作

3. Looking for Jobs  
好找工作 !!

## Information Versus Available Storage



Source: <http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>

Source: IDC, 2007

- Computational Load
- Genome Data  
**8x Growth / 18 month**
- Moore's Law  
**2x Growth / 18 months**

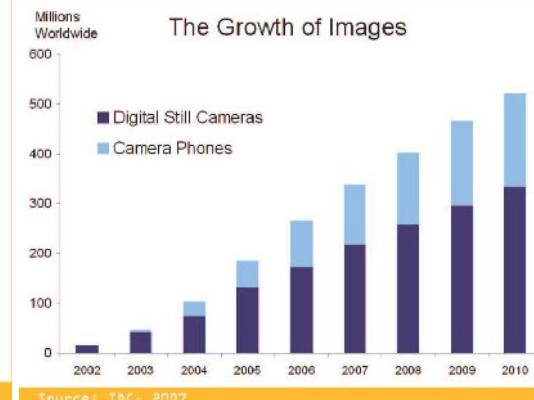
x Multiplier

## 2007 Data Explore

**Top 1 : Human Genomics – 7000 PB / Year**  
**Top 2 : Digital Photos – 1000 PB+/ Year**  
**Top 3 : E-mail (no Spam) – 300 PB+ / Year**



Source: IDC, 2007



Source: IDC, 2007



Total digital data to be created this year **270,000PB** (IDC)

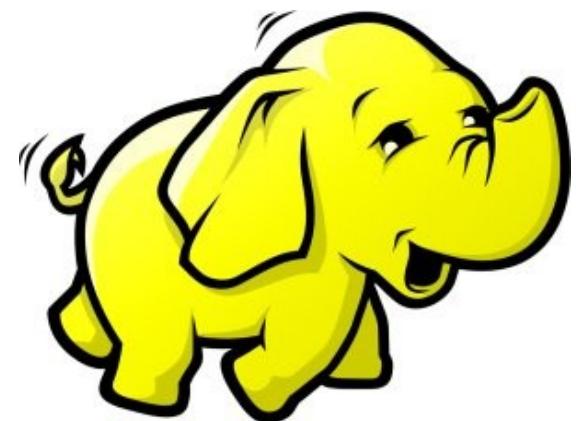
Source: [http://lib.stanford.edu/files/see\\_pasig\\_dic.pdf](http://lib.stanford.edu/files/see_pasig_dic.pdf)



# Hadoop 專業術語

## Introduction to Hadoop Terminology

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Two Key Elements of Operating System

## 作業系統兩大關鍵組成元素

Scheduler  
程序排程



File System  
檔案系統



# Terminologies of Hadoop

## *Hadoop* 文件中的專業術語

- Job
  - 任務
- Task
  - 小工作
- JobTracker
  - 任務分派者
- TaskTracker
  - 小工作的執行者
- Client
  - 發起任務的客戶端
- Map
  - 應對
- Reduce
  - 總和



- Namenode
  - 名稱節點
- Datanode
  - 資料節點
- Namespace
  - 名稱空間
- Replication
  - 副本
- Blocks
  - 檔案區塊 (64M)
- Metadata
  - 屬性資料



# Two Key Roles of HDFS

## HDFS 軟體架構的兩種關鍵角色

名稱節點

NameNode

- Master Node
- Manage NameSpace of HDFS
- Control Permission of Read and Write
- Define the policy of Replication
- Audit and Record the NameSpace
- Single Point of Failure

資料節點

DataNode

- Worker Nodes
- Perform operation of Read and Write
- Execute the request of Replication
- Multiple Nodes

# Two Key Roles of Job Scheduler

## 程序排程的兩種關鍵角色

### JobTracker

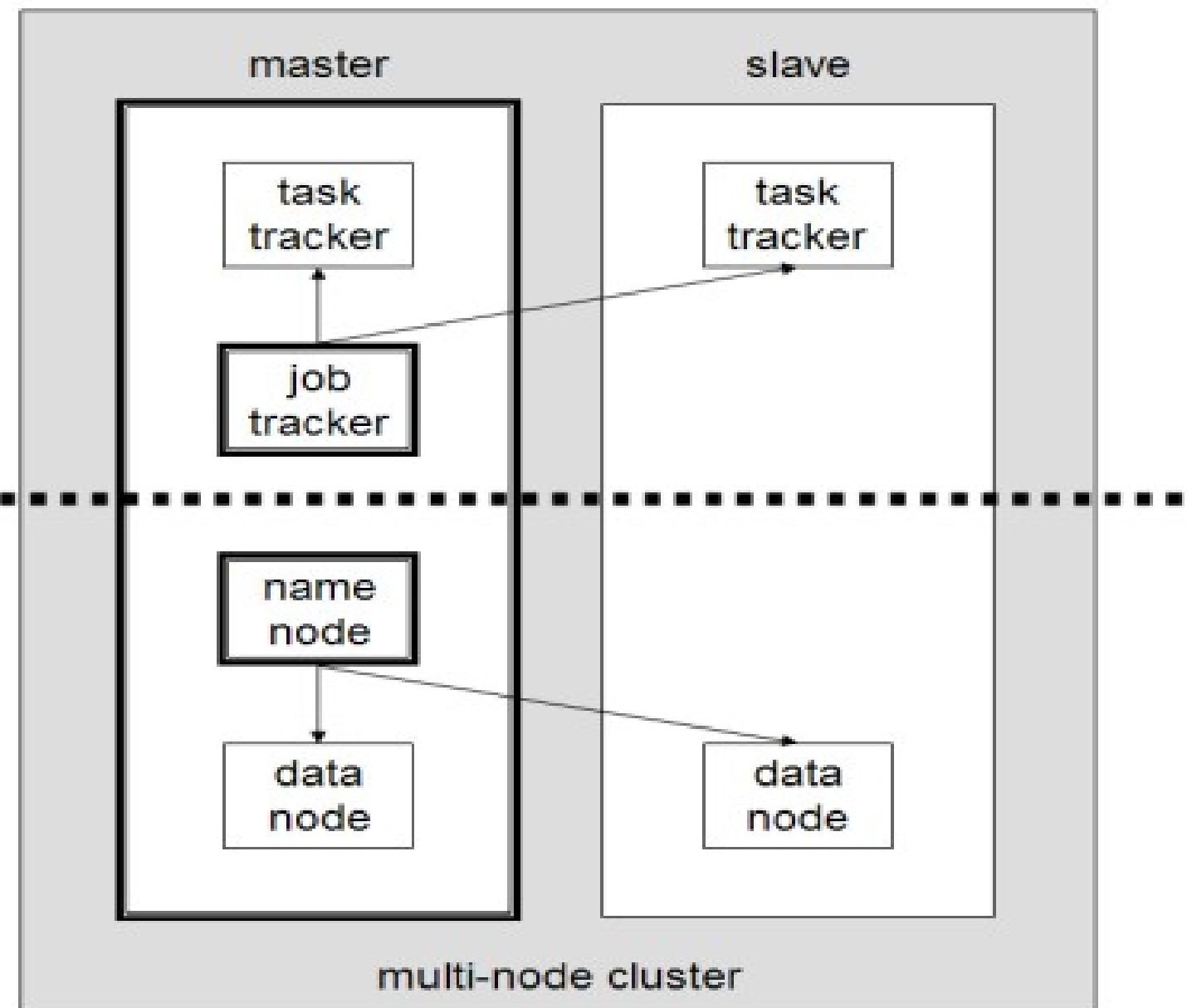
- Master Node
- Receive Jobs from Hadoop Clients
- Assigned Tasks to TaskTrackers
- Define Job Queuing Policy, Priority and Error Handling
- Single Point of Failure

### TaskTracker

- Worker Nodes
- Execute Mapper and Reducer Tasks
- Save Results and report task status
- Multiple Nodes

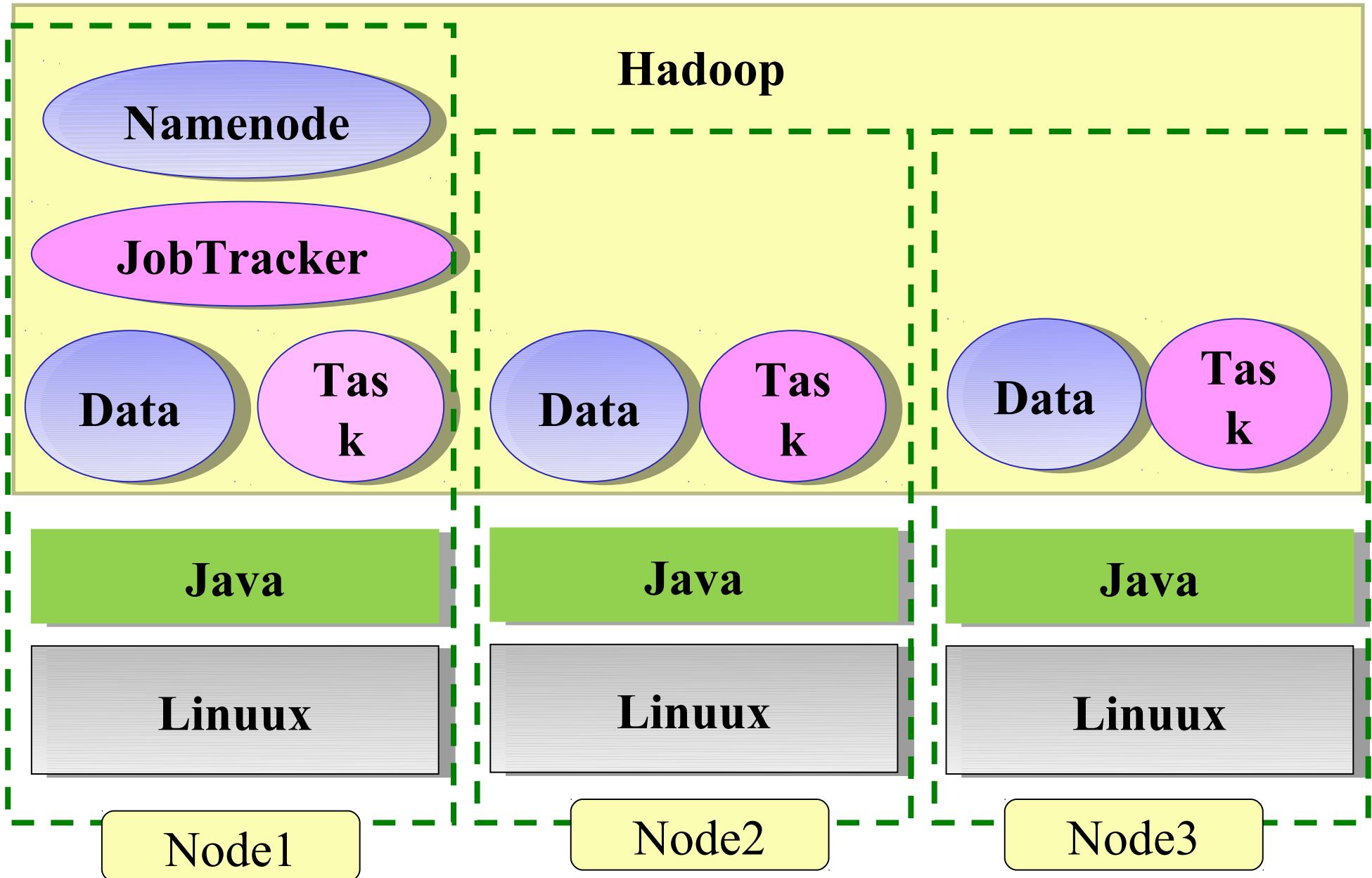
# Different Roles of Hadoop Architecture

## *Hadoop* 軟體架構中的不同角色



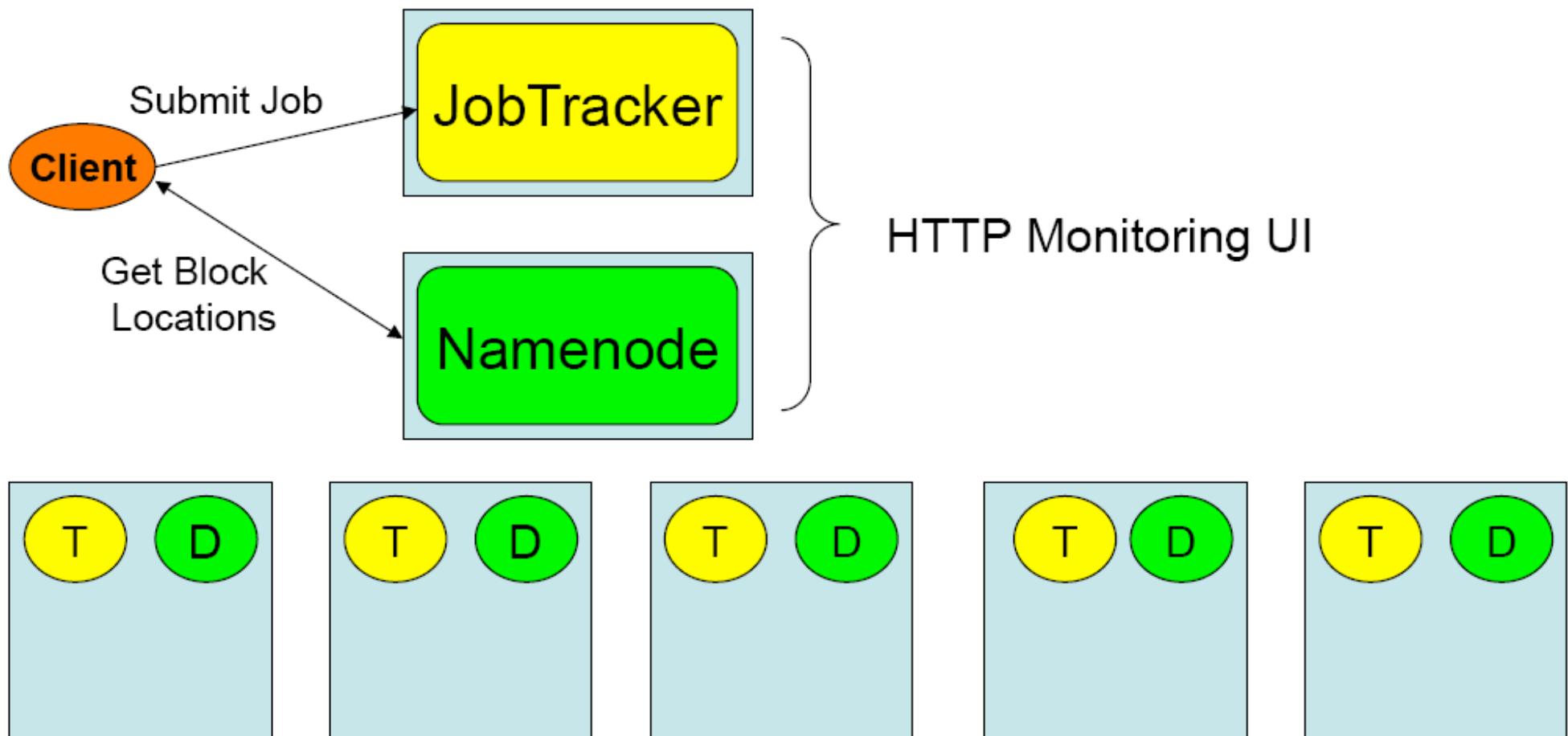
# Distributed Operating System of Hadoop

## Hadoop 建構成一個分散式作業系統



# About Hadoop Client ...

## 不在雲裡的 *Hadoop Client*



# What we learn today ?

WHAT

Hadoop 是運算海量資料的軟體平台 !!

hadoop is a software platform to process vast amount of data !!

WHO

始祖是 Doug Cutting , Apache 社群支持 , Yahoo 贊助

From Doug Cutting to Apache Community, Yahoo and more !

WHEN

Hadoop 是 2004 年從 Nutch 分裂出來的專案 !!

Hadoop became separate project since year 2004 !!

WHY

資料大爆炸、資料探勘、找工作

Data Explore, Data Mining, Jobs !!

HOW

建構在大型的個人電腦叢集之上

Install on large clusters built of commodity hardware !!



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



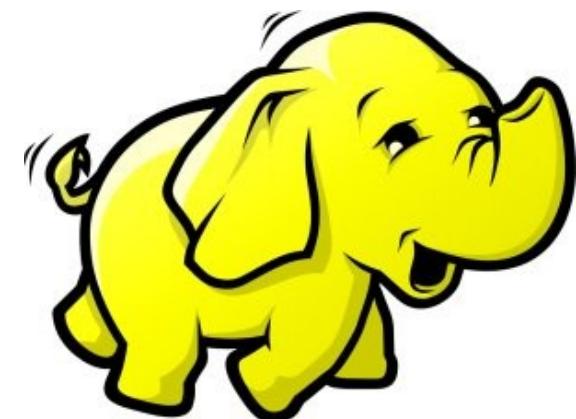
Powered by **DRBL**



# HDFS 簡介

Introduction to Hadoop Distributed File System

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# What is HDFS ??

## 什麼是 **HDFS** ??

- **Hadoop Distributed File System**

- 實現類似 Google File System 分散式檔案系統
- Reference from Google File System.
- 一個易於擴充的分散式檔案系統，目的為對大量資料進行分析
- A scalable distributed file system for large data analysis .
- 運作於廉價的普通硬體上，又可以提供容錯功能
- based on commodity hardware with high fault-tolerant.
- 紿大量的用戶提供總體性能較高的服務
- It have better overall performance to serve large amount of users.

# Features of HDFS ...

## HDFS 的特色是 ...

- 硬體錯誤容忍能力 **Fault Tolerance**
  - 硬體錯誤是正常而非異常
  - Failure is the norm rather than exception
  - 自動恢復或故障排除
  - automatic recovery or report failure
- 串流式的資料存取 **Streaming data access**
  - 批次處理多於用戶交互處理
  - Batch processing rather than interactive user access.
  - 高 Throughput 而非低 Latency
  - High aggregate data bandwidth (throughput)

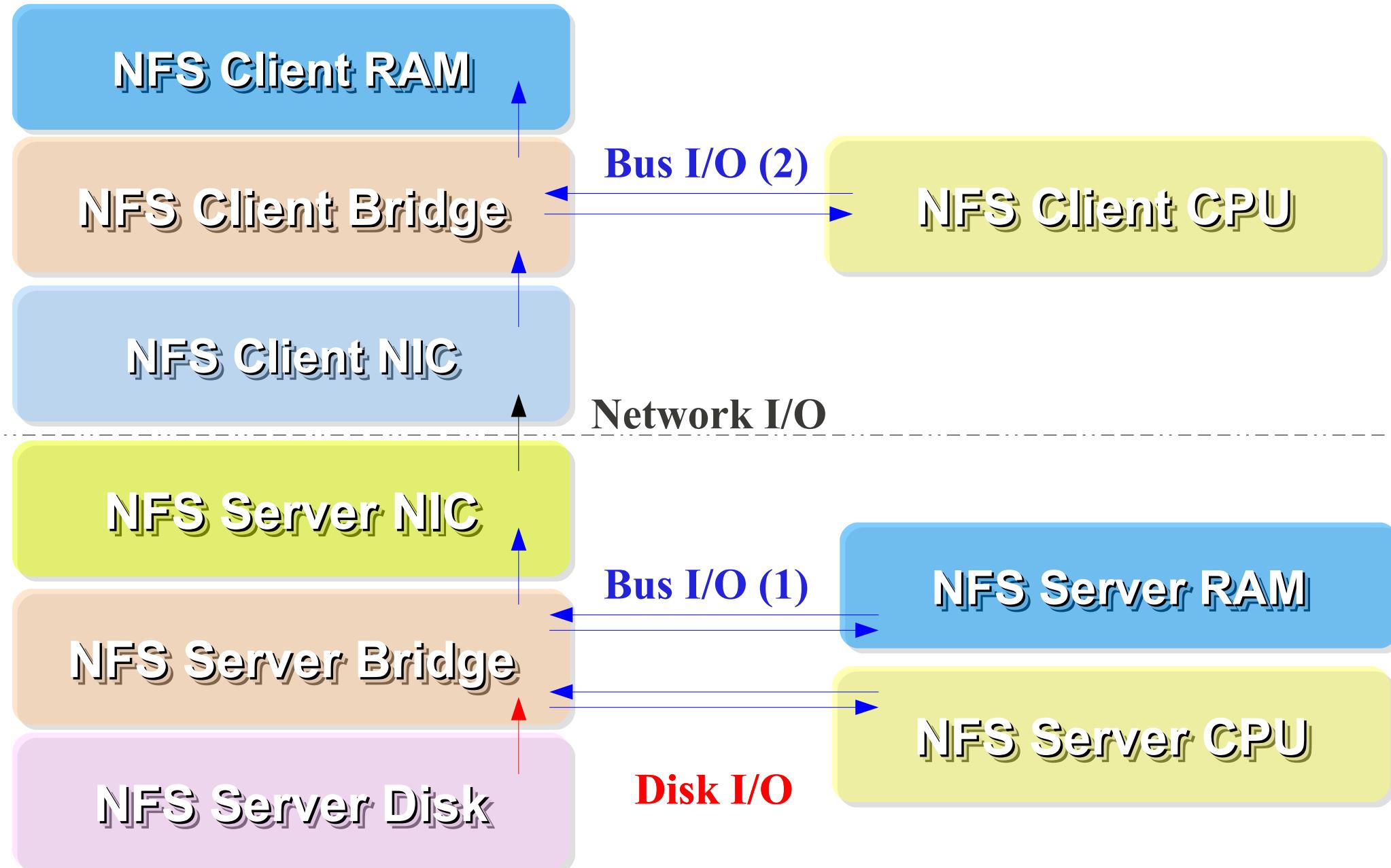
# Features of HDFS ...

## HDFS 的特色是 ...

- 大規模資料集 **Large data sets and files**
  - 支援 Petabytes 等級的磁碟空間
  - Support Petabytes size
- 一致性模型 **Coherency Model**
  - 一次寫入，多次存取 **Write-once-read-many**
  - 簡化一致性處理問題 **This assumption simplifies coherency**
- 在地運算 **Data Locality**
  - 到資料的節點上計算 > 將資料從遠端複製過來計算
  - “move compute to data” > “move data to compute”
- 異質平台移植性 **Heterogeneous**
  - 即使硬體不同也可移植、擴充
  - HDFS could be deployed on different hardware

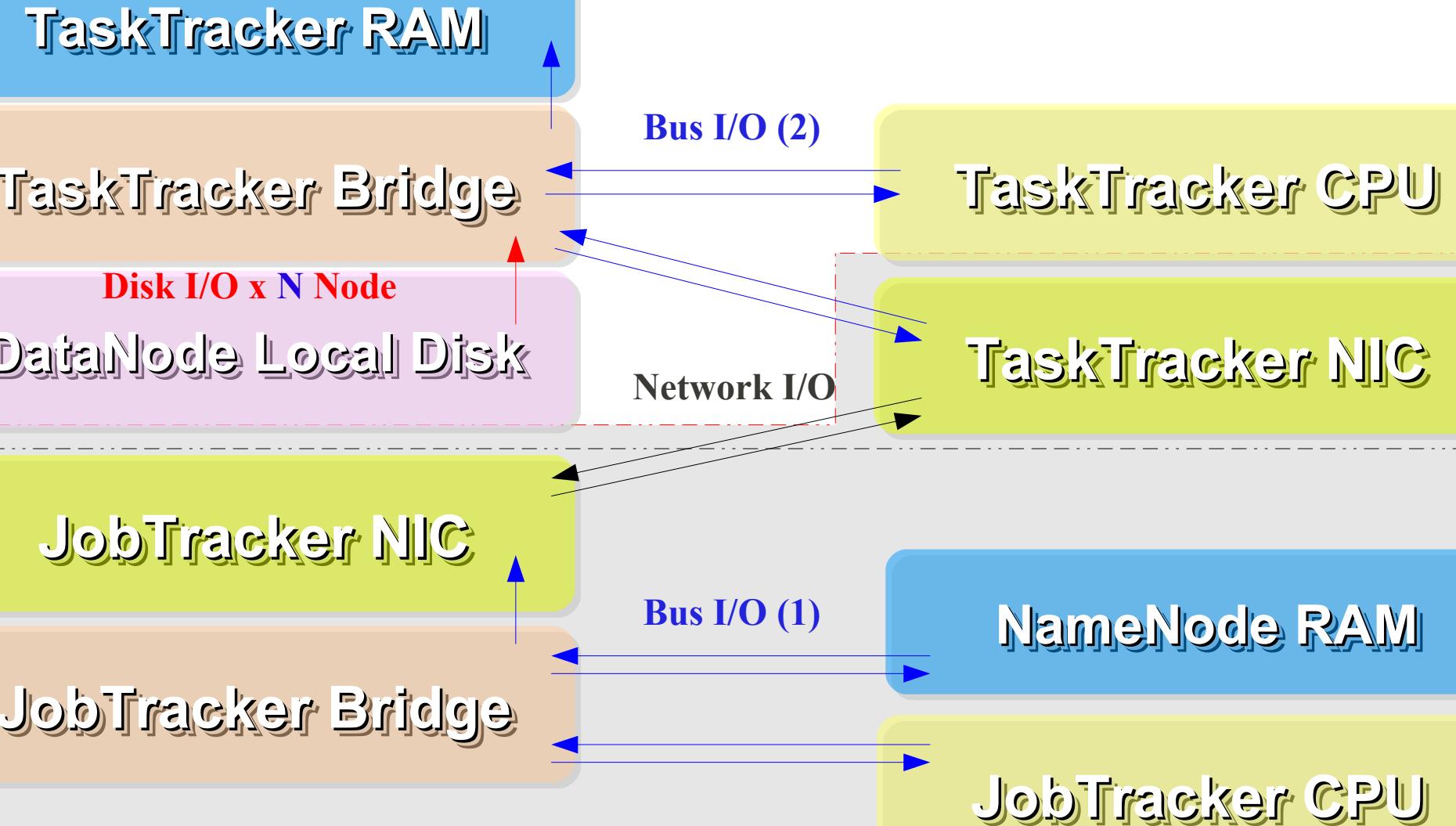
# Parallel Computing using NFS storage

使用 **NFS** 進行平行運算



# Parallel Computing using HDFS

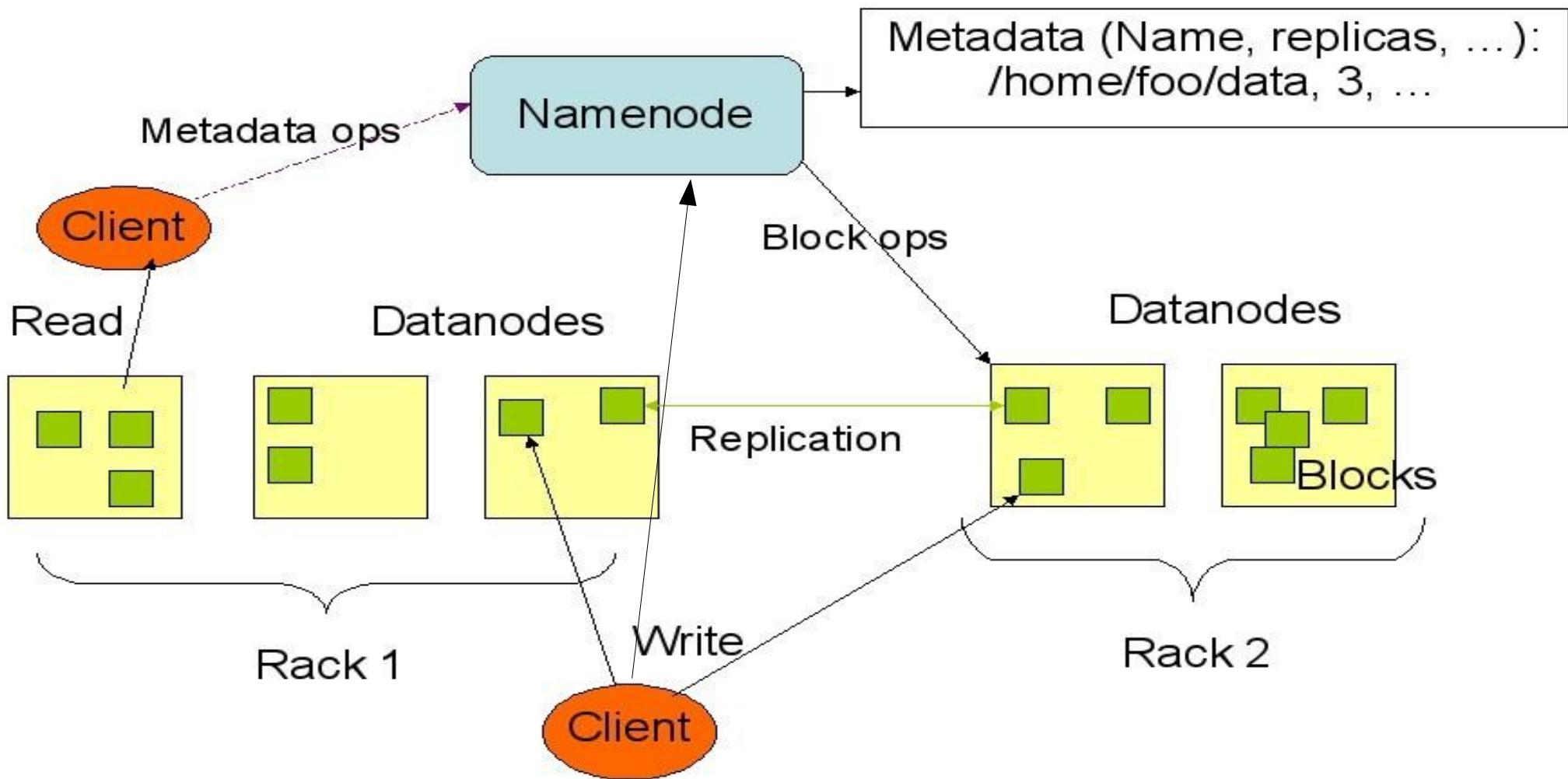
使用 **HDFS** 進行平行運算



# How HDFS manage data ...

**HDFS** 如何管理資料 ...

HDFS Architecture



# How does HDFS work ...

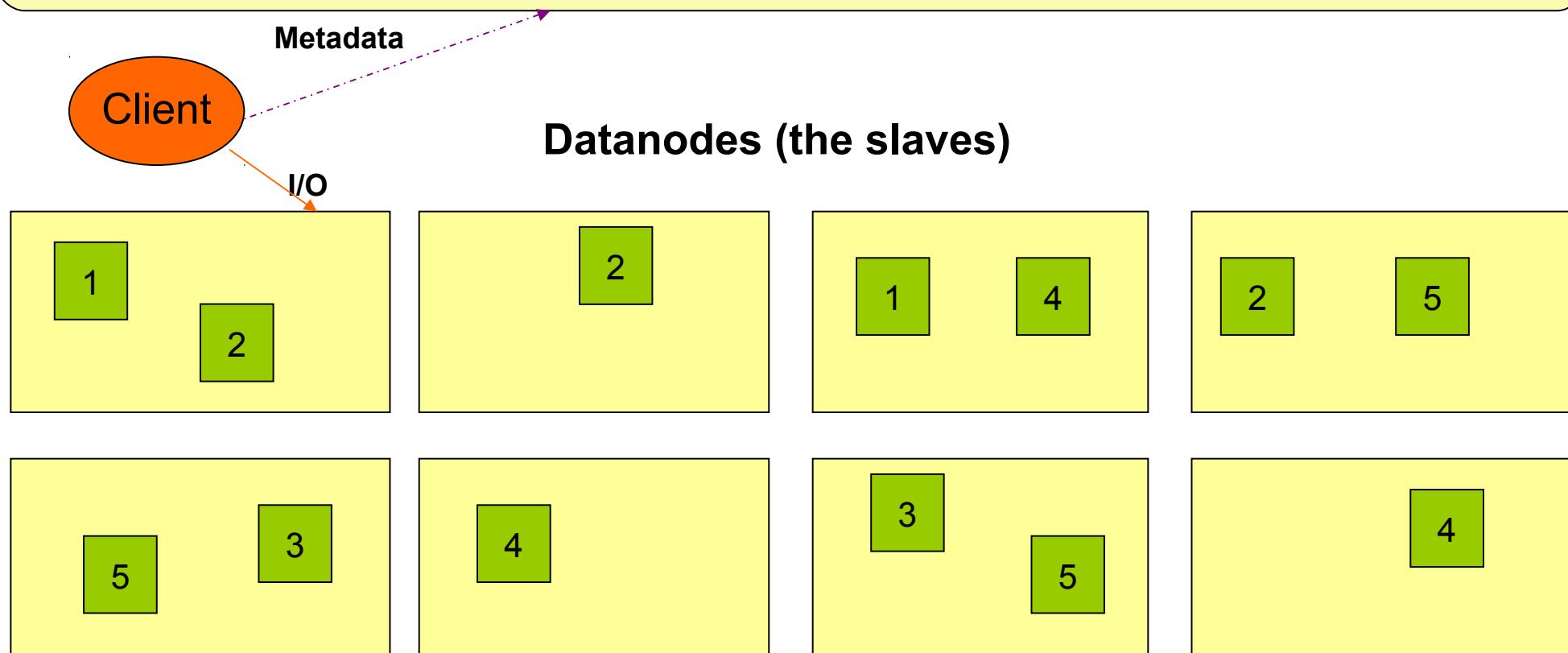
## HDFS 如何運作 ...

Namenode (the master)

Path and Filename – Replication , blocks

name:/users/joeYahoo/myFile - copies:2, blocks:{1,3}

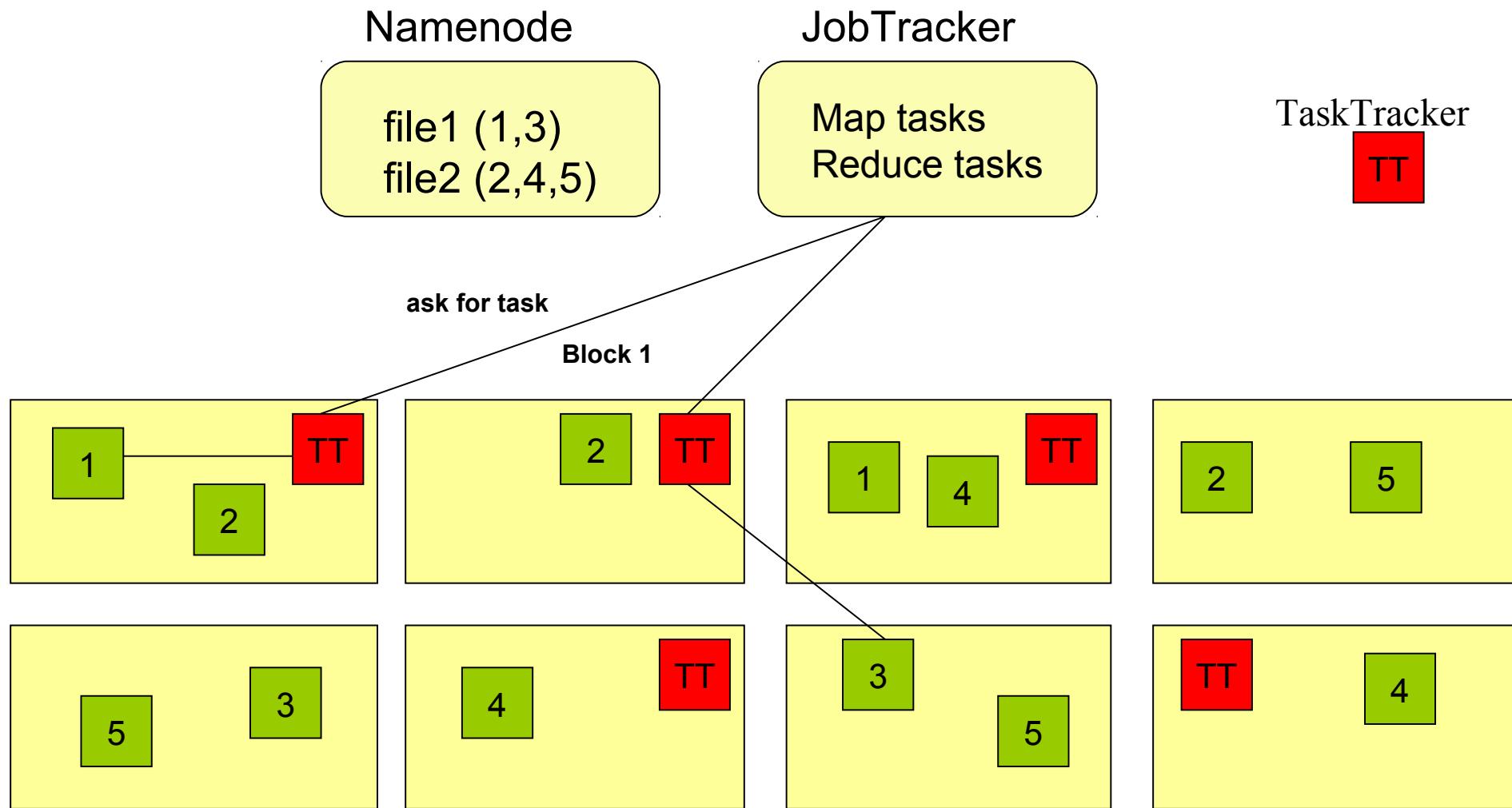
name:/users/bobYahoo/someData.gzip, copies:3, blocks:{2,4,5}



# About Data Locality ...

## HDFS 如何達成在地運算 ...

- Increase reliability and read bandwidth
  - robustness : read replication while found any failure
  - High read bandwith : distribute read ( but increase write bottleneck )



# About Fault Tolerance ...

## HDFS 如何達成容錯機制 ...

資料崩毀  
Data Corrupt

網路或資料  
節點失效  
Network Fault  
DataNode Fault

名稱節點錯誤  
NameNode Fault

- 資料完整性 Data integrity
  - checked with CRC32
  - 用副本取代出錯資料
  - Replace corrupt block with replication one
- Heartbeat
  - Datanode send **heartbeat** to Namenode
- Metadata
  - FSImage 、 Editlog 為核心印象檔及日誌檔
  - FSImage – core file system mapping image
  - Editlog – like. SQL transaction log
  - 多份儲存，當名稱節點故障時可以手動復原
  - Multiple backups of FSImage and Editlog
  - Manually recovery while NameNode Fault

# Coherency Model and Performance of HDFS

## HDFS 的一致性機制與效能 ...

- 檔案一致性機制 **Coherency model of files**
  - 刪除檔案＼新增寫入檔案＼讀取檔案皆由名稱節點負責
  - NameNode handle the operation of write, read and delete.
- 巨量空間及效能機制 **Large Data Set and Performance**
  - 預設每個區塊大小以 64MB 為單位
  - By default, the block size is 64MB
  - 大區塊可提高存取效率
  - Bigger block size will enhance read performance
  - 檔案有可能大過一顆磁碟
  - Single file stored on HDFS might be larger than single physical disk of DataNode.
  - 區塊均勻散佈各節點以分散讀取流量
  - Fully distributed blocks increase throughput of reading.

# POSIX like HDFS commands

與 **POSIX** 相似的操作指令 ...

```
jazz@hadoop:~$ hadoop fs
Usage: java FsShell
      [-ls <path>]
      [-lsr <path>]
      [-du <path>]
      [-dus <path>]
      [-count[-q] <path>]
      [-mv <src> <dst>]
      [-cp <src> <dst>]
      [-rm <path>]
      [-rmr <path>]
      [-expunge]
      [-put <localsrc> ... <dst>]
      [-copyFromLocal <localsrc> ... <dst>]
      [-moveFromLocal <localsrc> ... <dst>]
      [-get [-ignoreCrc] [-crc] <src> <localdst>]
      [-getmerge <src> <localdst> [addnl]]
      [-cat <src>]
      [-text <src>]
      [-copyToLocal [-ignoreCrc] [-crc] <src> <localdst>]
      [-moveToLocal [-crc] <src> <localdst>]
      [-mkdir <path>]
      [-setrep [-R] [-w] <rep> <path/file>]
      [-touchz <path>]
      [-test -[ezd] <path>]
      [-stat [format] <path>]
      [-tail [-f] <file>]
      [-chmod [-R] <MODE[,MODE]... | OCTALMODE> PATH...]
      [-chown [-R] [OWNER][:[GROUP]] PATH...]
      [-chgrp [-R] GROUP PATH...]
      [-help [cmd]]
```



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



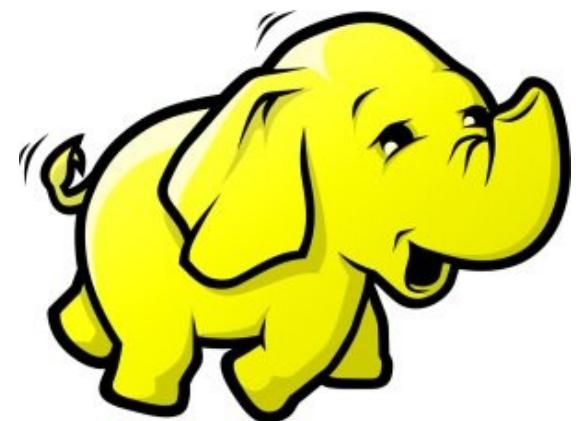
Powered by **DRBL**



# MapReduce 簡介

## Introduction to MapReduce

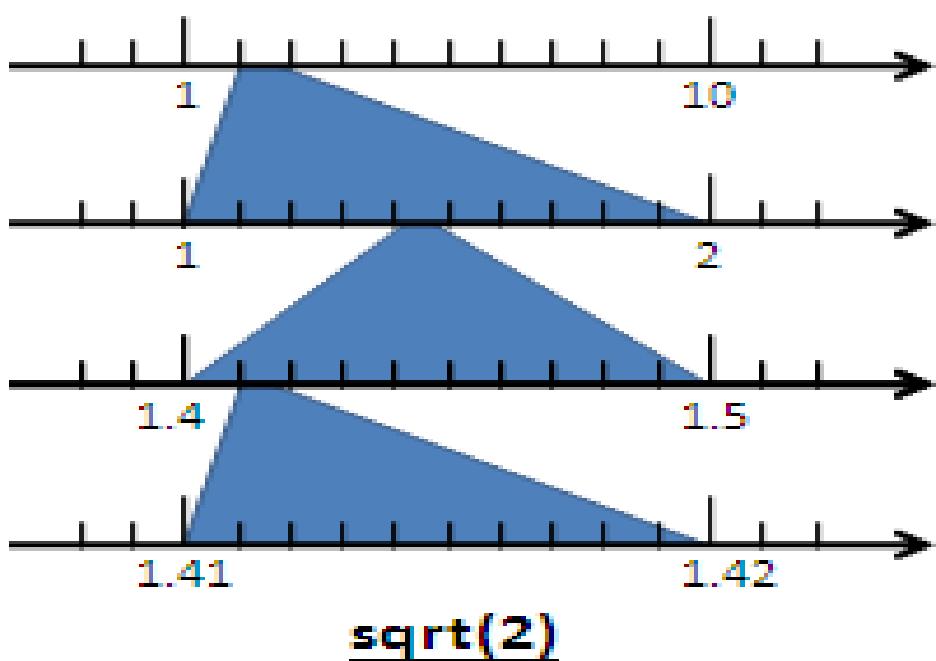
**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



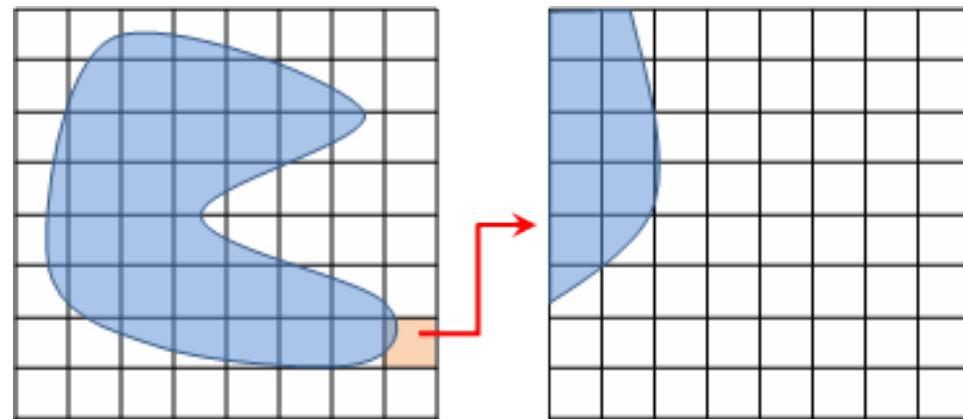
# Divide and Conquer Algorithms

分而治之演算法

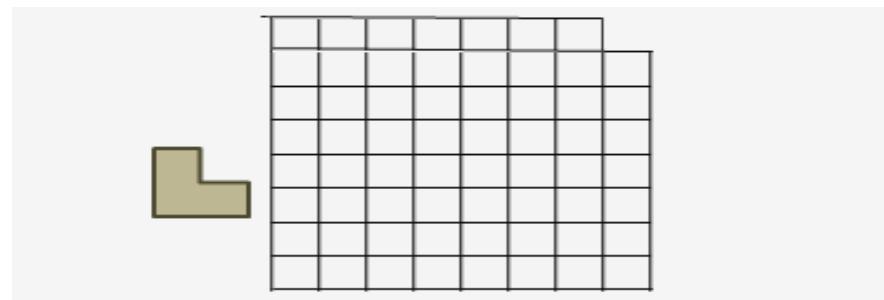
Example 1:



Example 2:



Example 3:



Example 4: The way to climb 5 steps stair within 2 steps each time. 眼前有五階樓梯，每次可踏上一階或踏上兩階，那麼爬完五階共有幾種踏法？

Ex : (1,1,1,1,1) or (1,2,1,1)

# What is MapReduce ??

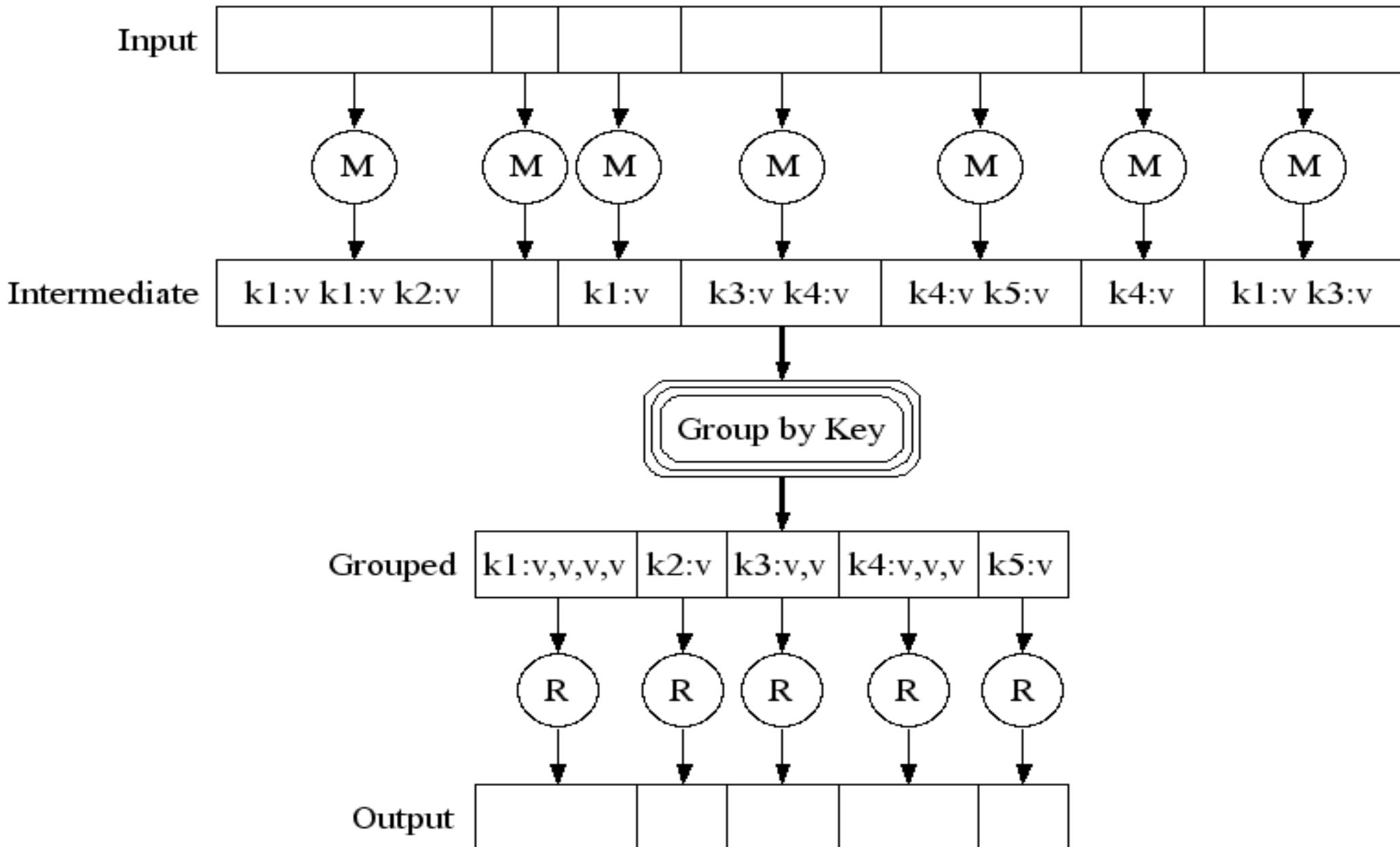
## 什麼是 **MapReduce** ??

- MapReduce 是 Google 申請的軟體專利，主要用來處理大量資料
- **MapReduce is a patented software framework introduced by Google to support distributed computing on large data sets on clusters of computers.**
- 啓發自函數編程中常用的 map 與 reduce 函數。
- **The framework is inspired by map and reduce functions commonly used in functional programming, although their purpose in the MapReduce framework is not the same as their original forms**
  - Map(...):  $N \rightarrow N$ 
    - Ex.  $[ 1,2,3,4 ] - (*2) \rightarrow [ 2,4,6,8 ]$
  - Reduce(...):  $N \rightarrow 1$ 
    - $[ 1,2,3,4 ] - (\text{sum}) \rightarrow 10$
- **Logical view of MapReduce**
  - $\text{Map}(k1, v1) \rightarrow \text{list}(k2, v2)$
  - $\text{Reduce}(k2, \text{list } (v2)) \rightarrow \text{list}(k3, v3)$

Source: <http://en.wikipedia.org/wiki/MapReduce>

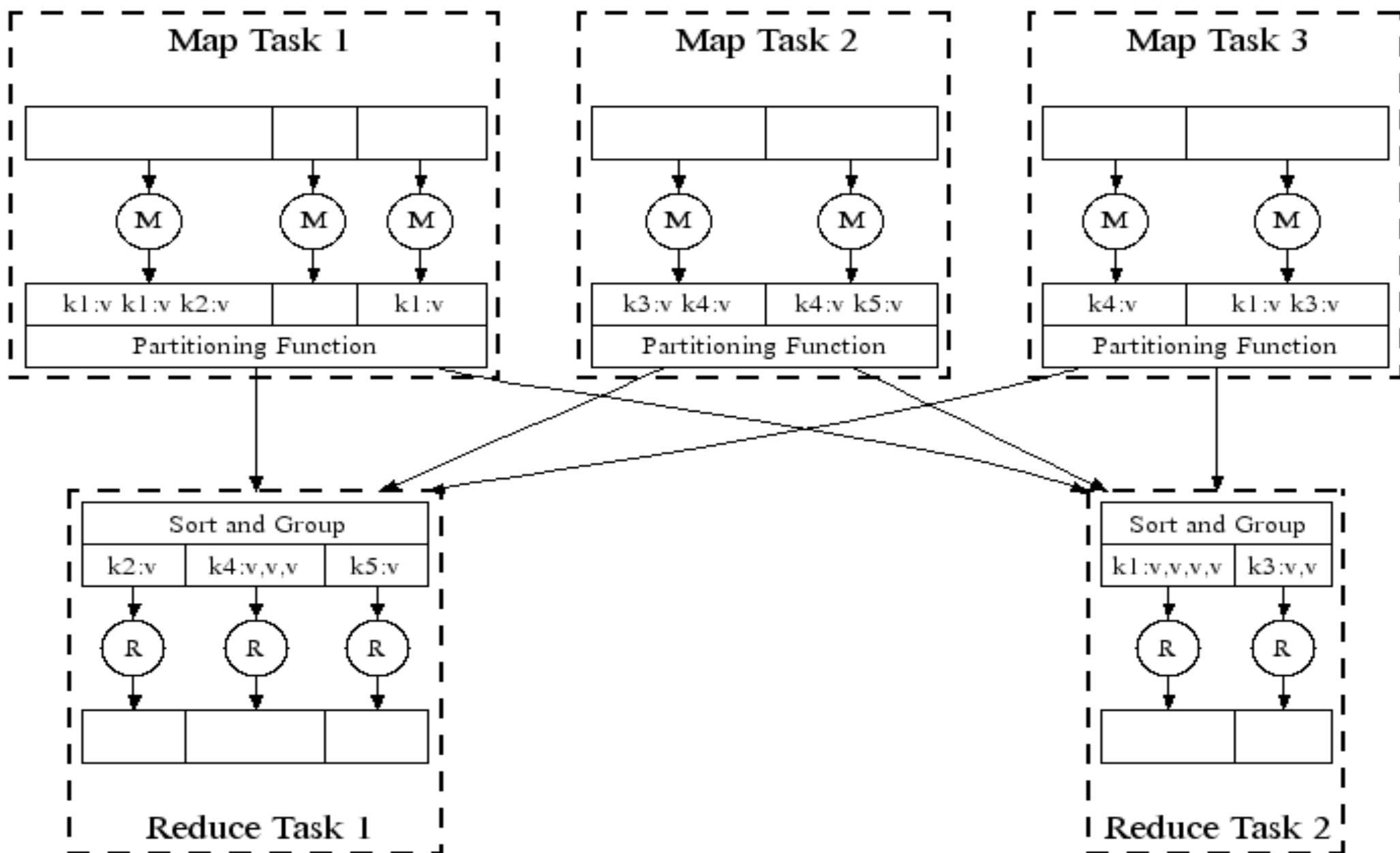
# Google's MapReduce Diagram

## Google 的 *MapReduce* 圖解



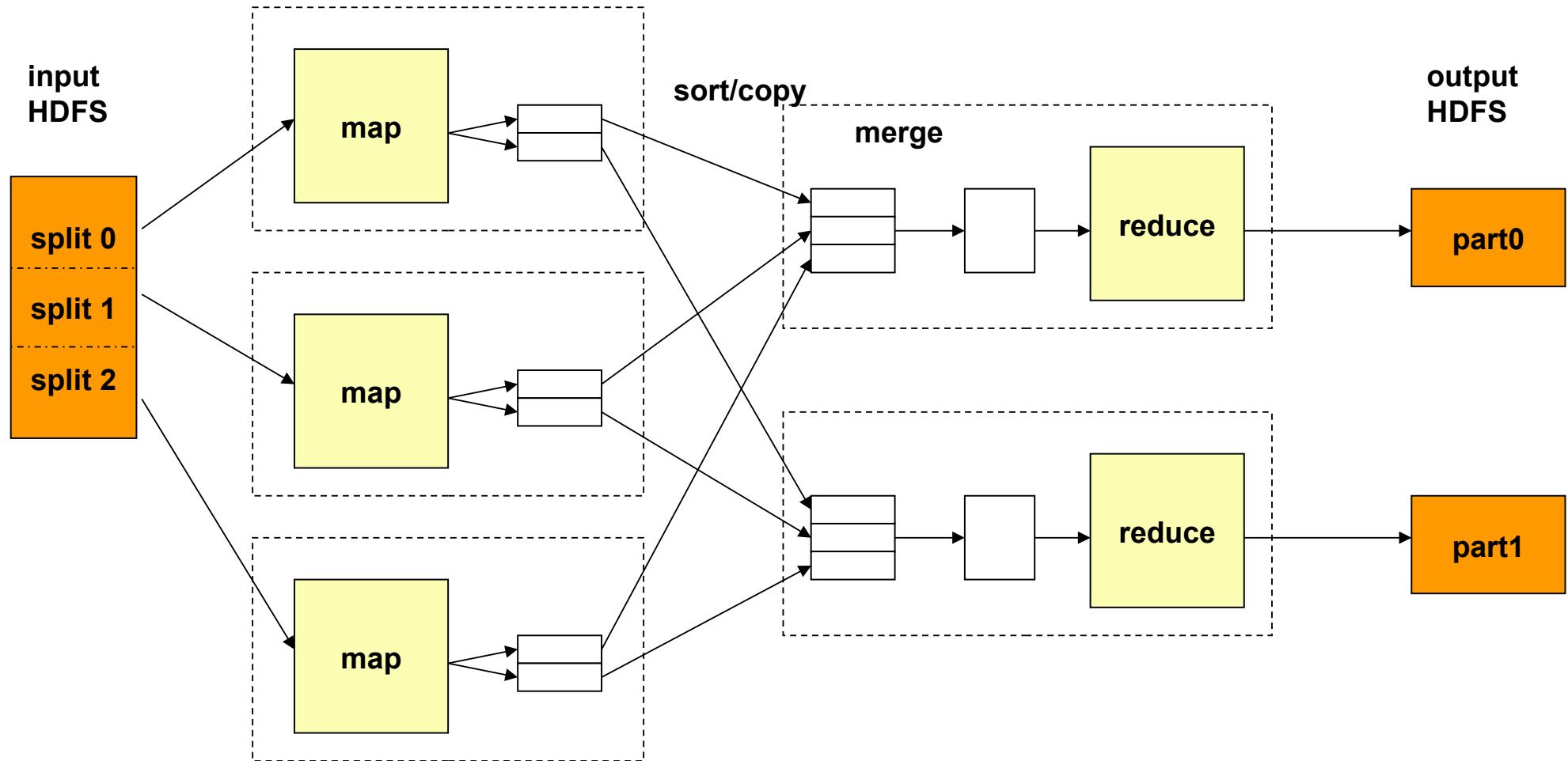
# Google's MapReduce in Parallel

## Google 的 *MapReduce* 平行版圖解



# How does MapReduce work in Hadoop

## Hadoop MapReduce 運作流程



JobTracker 跟 NameNode 取得需要運算的 blocks

JobTracker 選數個 TaskTracker 來作 Map 運算，產生些中間檔案

JobTracker 將中間檔案整合排序後，複製到需要的 TaskTracker 去

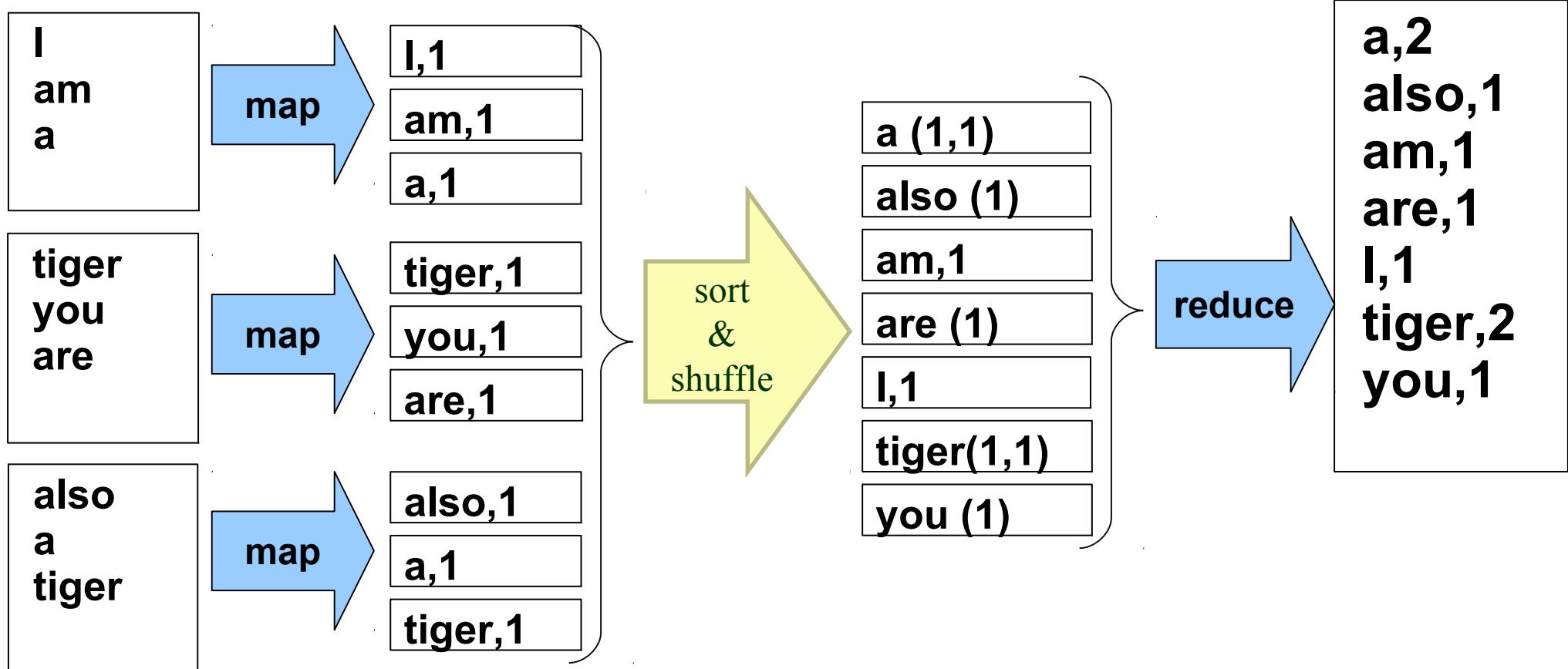
JobTracker 派遣 TaskTracker 作 reduce

reduce 完後通知 JobTracker 與 Namenode 以產生 output

# MapReduce by Example (1)

## MapReduce 運作實例 (1)

I am a tiger, you are also a tiger



JobTracker 先選了三個  
Tracker 做 map

Map 結束後，hadoop 進行  
中間資料的重組與排序

JobTracker 再選一個  
TaskTracker 作 reduce

# MapReduce by Example (2)

## MapReduce 運作實例 (2)

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \rightarrow \begin{bmatrix} \sqrt{a+b} \\ \sqrt{c+d} \end{bmatrix}$$

$$\begin{bmatrix} 1.0 & 0.0 & 3.0 \\ 3.2 & 0.8 & 32.0 \\ 1.0 & 14.0 & 1.0 \end{bmatrix} \rightarrow ?$$

(0,  $\sqrt{1.0 + 0.0 + 3.0}$ )  
(1,  $\sqrt{3.2 + 0.8 + 32.0}$ )  
(2,  $\sqrt{1.0 + 14.0 + 1.0}$ )

### Input File

```
0 0 1.0 // A[0][1] = 1.0
0 1 0.0 // A[0][1] = 0.0
0 2 3.0 // A[0][2] = 3.0
1 0 3.2 // A[1][0] = 3.2
1 1 0.8 // A[1][1] = 0.8
```

map

(0, 1.0)  
(0, 0.0)  
(0, 3.0)  
(1, 3.2)  
(1, 0.8)

reduce

```
1 2 32.0 // A[1][2] = 32.0
2 0 1.0 // A[2][0] = 1.0
2 1 14.0 // A[2][1] = 14.0
2 2 1.0 // A[2][2] = 1.0
```

map

(1, 32.0)  
(2, 1.0)  
(2, 14.0)  
(2, 1.0)

sort / merge

(0, {1.0, 0.0, 3.0})  
(1, {3.2, 0.8, 32.0})  
(2, {1.0, 14.0, 1.0})

# MapReduce is suitable to ....

## *MapReduce* 合適用於 ....

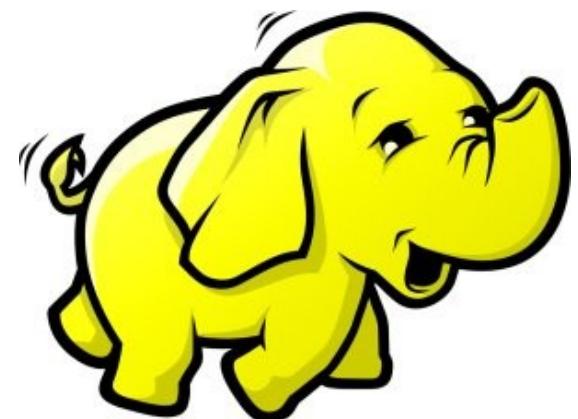
- 大規模資料集
  - Large Data Set
    - Text tokenization
    - Indexing and Search
    - Data mining
    - machine learning
    - ...
  - 可拆解
  - Parallelization
- 
- <http://www.dbms2.com/2008/08/26/known-applications-of-mapreduce/>
  - <http://wiki.apache.org/hadoop/PoweredBy>



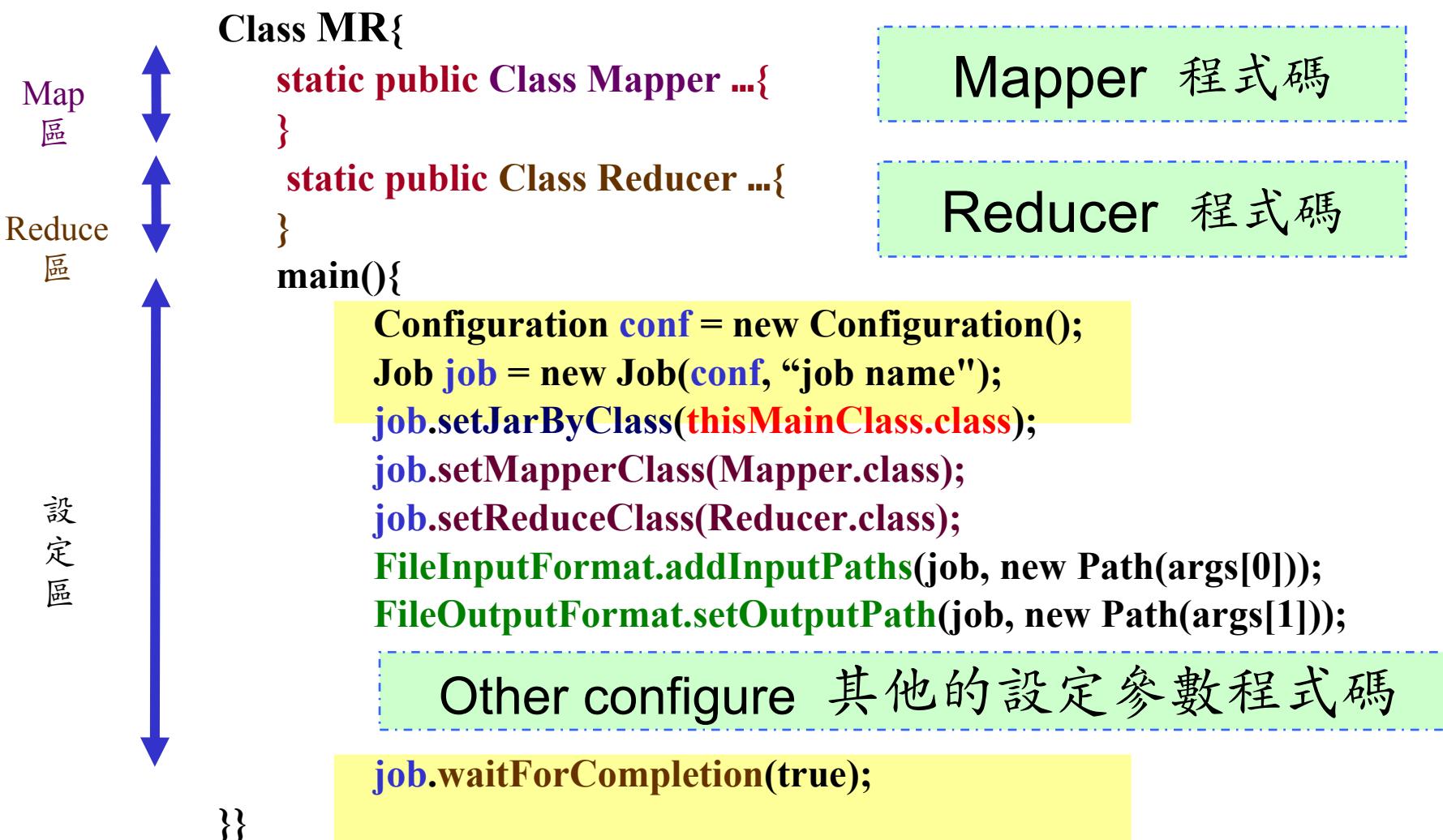
# MapReduce 程式設計入門

## MapReduce Programming 101

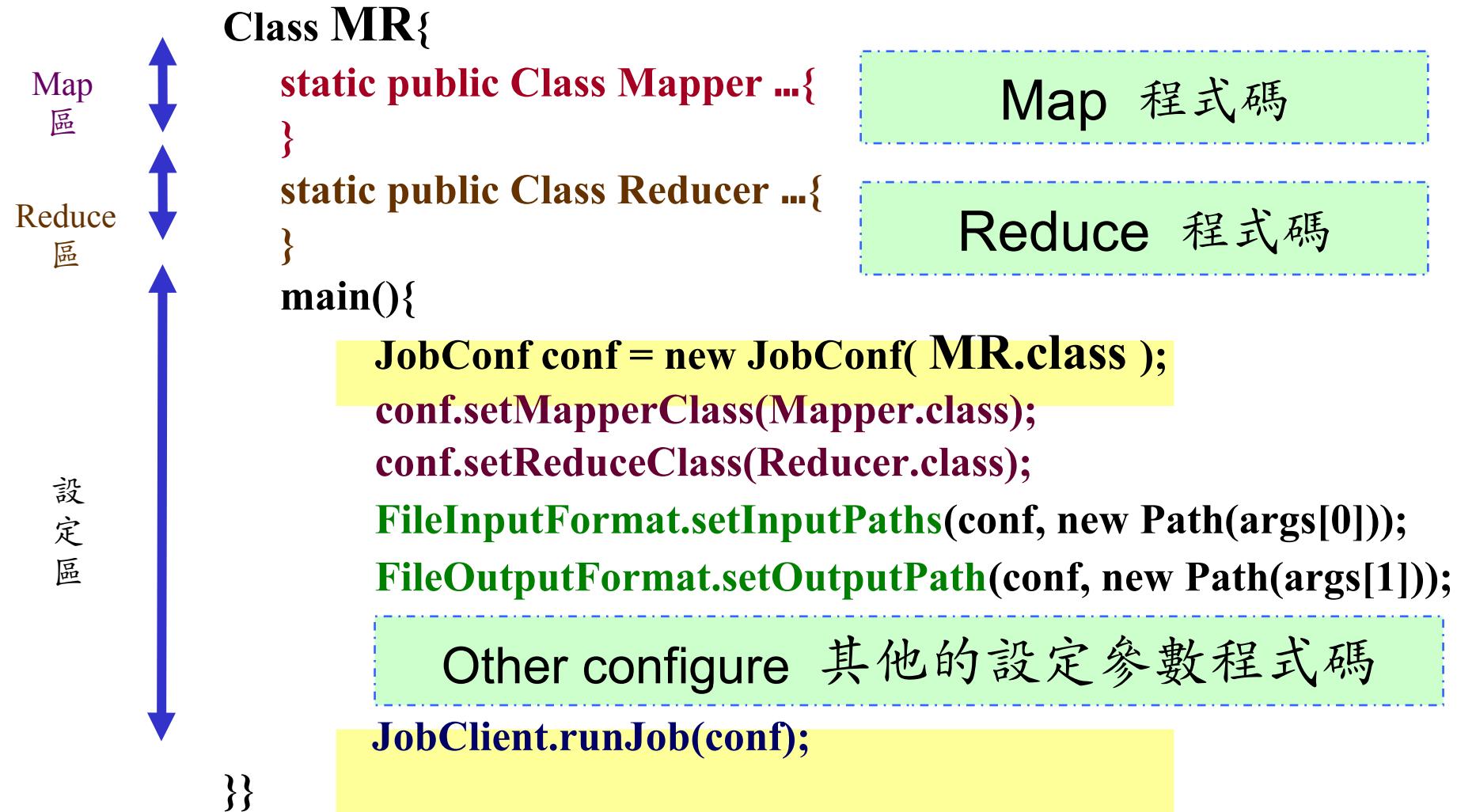
**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Program Prototype (v 0.20)

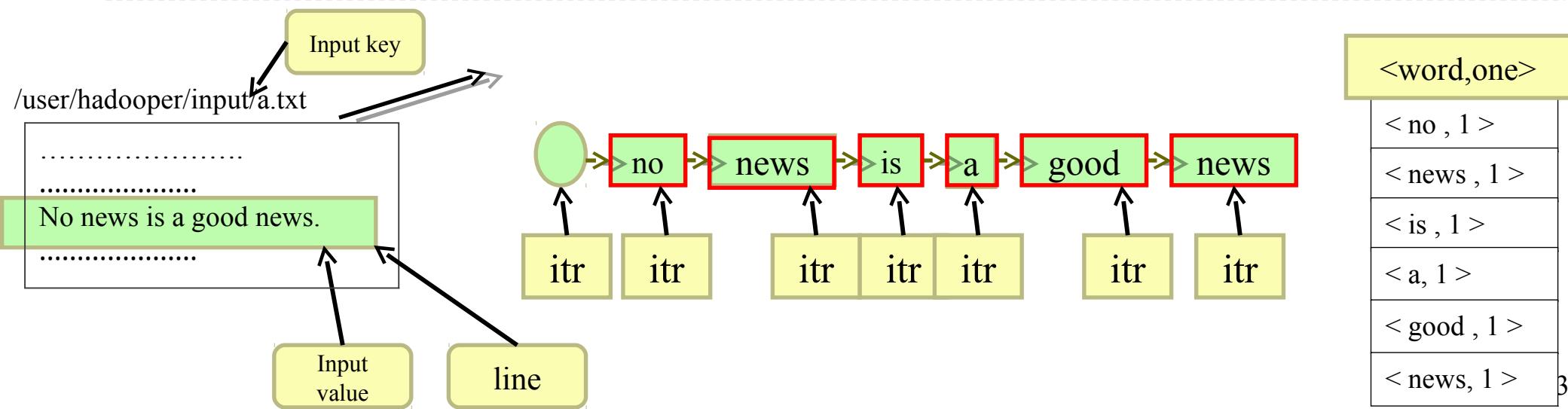


# Program Prototype (v 0.18)



# Word Count - mapper

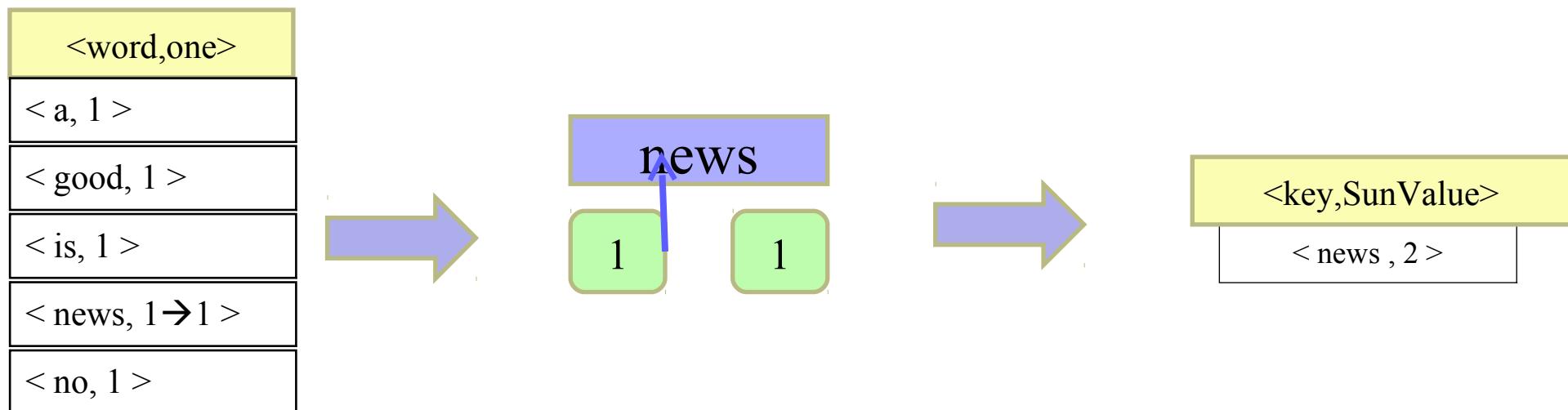
```
1 class MyMapper extends Mapper<LongWritable, Text, Text, IntWritable> {  
2     private final static IntWritable one = new IntWritable(1);  
3     private Text word = new Text();  
4     public void map( LongWritable key, Text value, Context context)  
5         throws IOException , InterruptedException {  
6         String line = ((Text) value).toString();  
7         StringTokenizer itr = new StringTokenizer(line);  
8         while (itr.hasMoreTokens()) {  
9             word.set(itr.nextToken());  
10            context.write(word, one);  
11        }  
12    }  
13}
```



# Word Count - reducer

```
1 class MyReducer extends Reducer< Text, IntWritable, Text, IntWritable> {  
2     IntWritable result = new IntWritable();  
3     public void reduce( Text key, Iterable <IntWritable> values, Context context)  
4         throws IOException, InterruptedException {  
5             int sum = 0;  
6             for( IntWritable val : values ) {  
7                 sum += val.get();  
8             }  
9             result.set(sum);  
10            context.write ( key, result);  
11        }  
12    }
```

~~for ( int i ; i < values.length ; i ++ ){  
 sum += values[i].get()  
}~~



# Word Count – main program

```
Class WordCount{  
    main()  
        Configuration conf = new Configuration();  
        Job job = new Job(conf, "job name");  
        job.setJarByClass(thisMainClass.class);  
        job.setMapperClass(MyMapper.class);  
        job.setReduceClass(MyReducer.class);  
        FileInputFormat.addInputPaths(job, new Path(args[0]));  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
        job.waitForCompletion(true);  
    }}
```



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



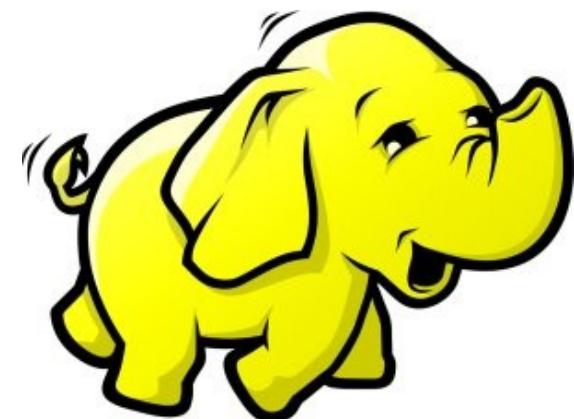
Powered by **DRBL**



# Hadoop 簡集設定解說

## Setup Hadoop Fully Distributed Mode

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Yahoo's Hadoop Cluster

## 雅虎的大象軍團

- ~10,000 machines running Hadoop in US
- The largest cluster is currently 2000 nodes
- Nearly 1 petabyte of user data (compressed, unreplicated)
- Running roughly 10,000 research jobs / week



# Hadoop Pseudo-Distributed Mode

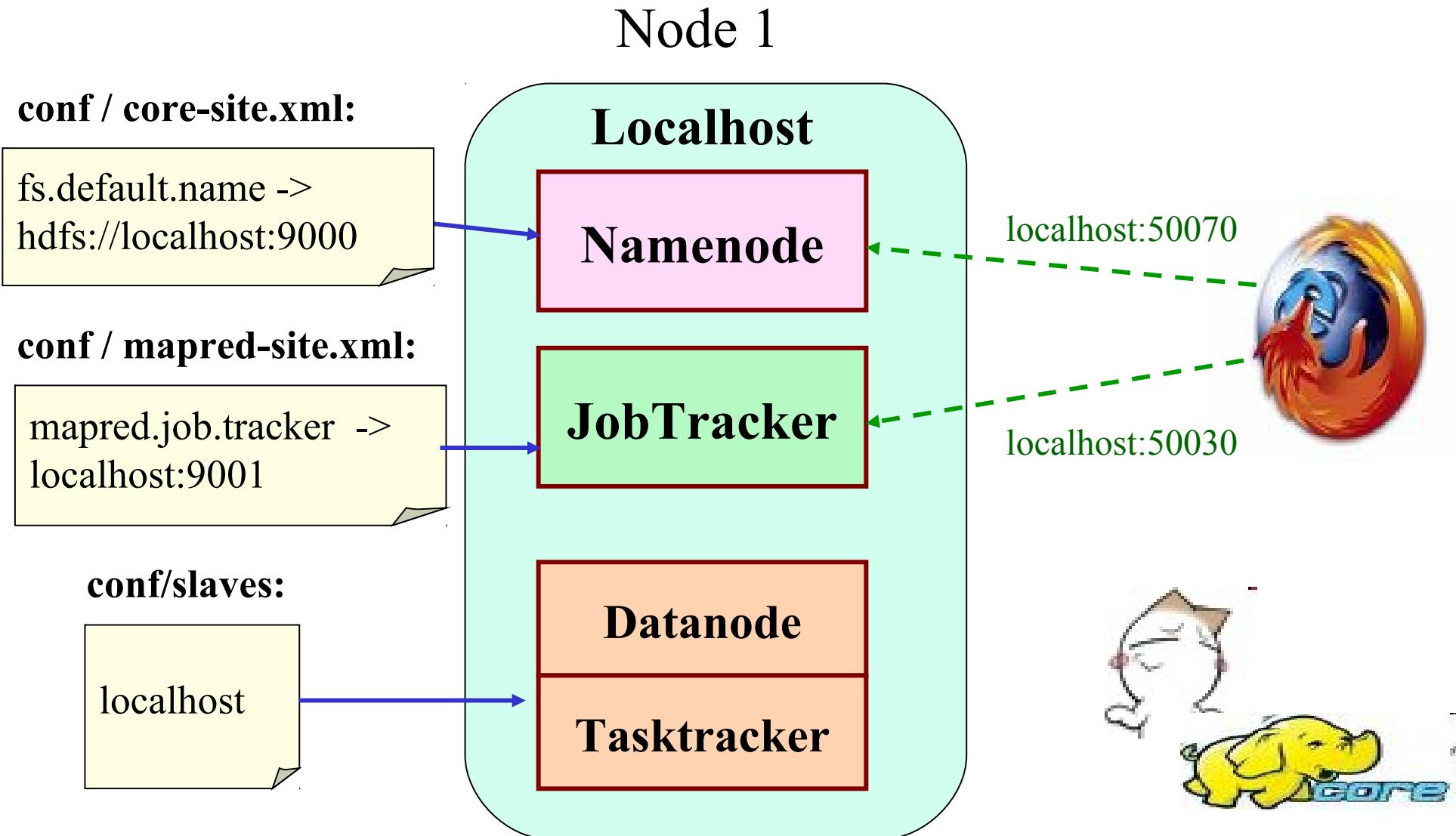
我們已經實作過單機模式

- Step 1: Setup SSH key exchange
- Step 2. Install Java
- Step 3: Download Hadoop Source Package
- Step 4: Configure hadoop-env.sh
  - `export JAVA_HOME=/usr/lib/jvm/java-6-sun`
- Step 5: Configure \*-site.xml
  - Set Namenode to `hdfs://localhost:9000`
  - Set Jobtracker to `localhost:9001`
  - `bin/hadoop namenode -format`
- Step 6: Format HDFS
- Step 7: Start Hadoop
  - `bin/start-all.sh`
- Step 8: Complete!! Let's check the status of Hadoop
  - Job admin <http://localhost:50030/> HDFS <http://localhost:50070/>



# Diagram of Pseudo-Distributed Mode

## Hadoop 單機環境示意圖



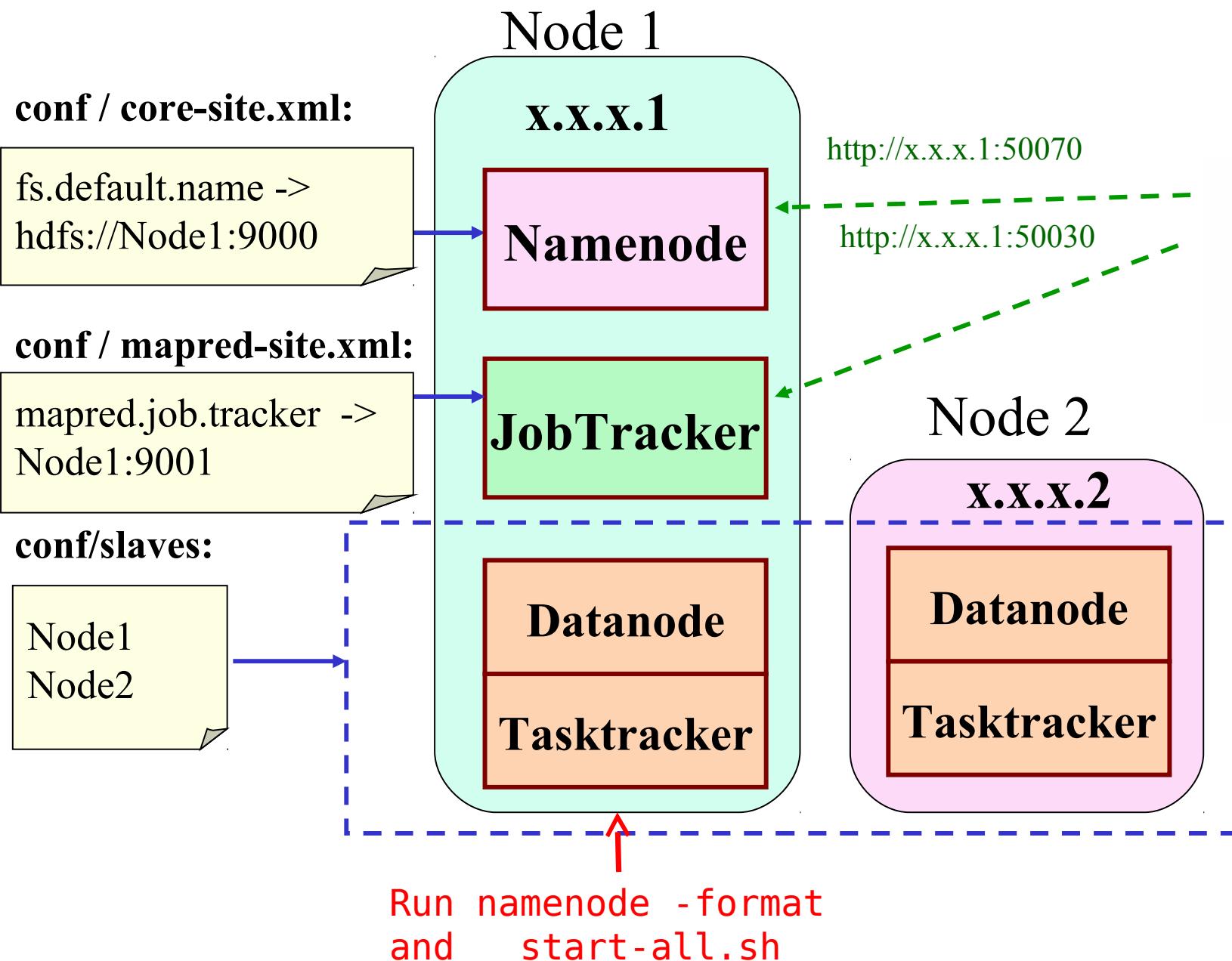
# Hadoop Fully-Distributed Mode

## 我們接著要用兩台電腦實作叢集模式

- Step 1: Setup SSH key exchange
- Step 2. Install Java
- Step 3: Download Hadoop Source Package
- Step 4: Configure `hadoop-env.sh`
  - `export JAVA_HOME=/usr/lib/jvm/java-6-sun`
- Step 5: Configure `*-site.xml`
  - Set Namenode to `hdfs://x.x.x.1:9000`
  - Set Jobtracker to `x.x.x.2:9001`
- Step 6: Configure Slaves
- Step 7: Syncronization of all slaves
- Step 8: Format HDFS
  - `bin/hadoop namenode -format`
- Step 9: Start Hadoop
  - On NameNode : `bin/start-dfs.sh`
  - On JobTracker : `bin/start-mapred.sh`
- Step 10: Complete!! Let's check the status of Hadoop
  - Job admin <http://x.x.x.2:50030/> HDFS <http://x.x.x.1:50070/>

# Use case #1

## 設定情境一



# Use case #2

## 設定情境二

conf / core-site.xml:

fs.default.name ->  
hdfs://Node1:9000

conf / mapred-site.xml:

mapred.job.tracker ->  
Node2:9001

conf/slaves:

Node1  
Node2

Node 1

Node 2

x.x.x.1

x.x.x.2

Namenode

JobTracker

Datanode

Datanode

Tasktracker

Tasktracker

http://x.x.x.1:50070



http://x.x.x.2:50030



Run namenode -format run start-mapred.sh  
and start-dfs.sh

# Use case #3

## 設定情境三

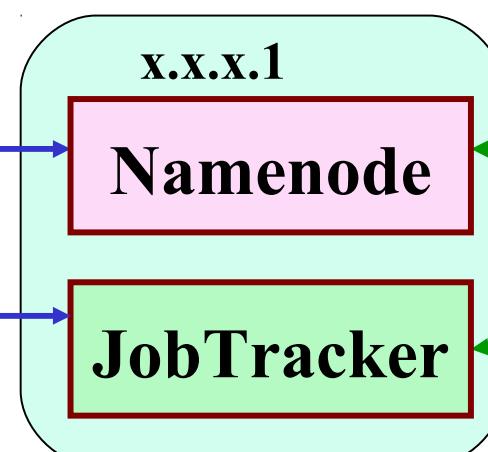
conf / core-site.xml:

fs.default.name ->  
hdfs://Node1:9000

conf / mapred-site.xml:

mapred.job.tracker ->  
Node1:9001

Node 1



http://x.x.x.1:50070

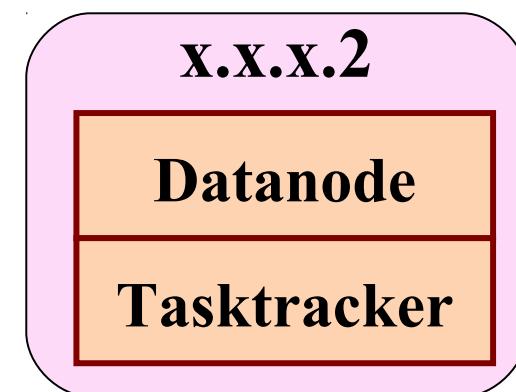
http://x.x.x.1:50030



conf/slaves:

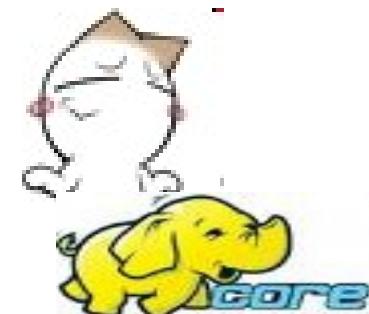
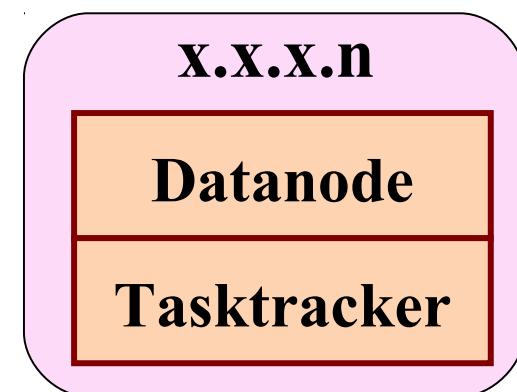
Node2  
.....  
NodeN

Node 2



...

Node N



# Use case #4

## 設定情境四

**conf / core-site.xml:**

fs.default.name ->  
hdfs://Node1:9000

Client



http://x.x.x.2:50030

**conf / mapred-site.xml:**

mapred.job.tracker ->  
Node2:9001

G



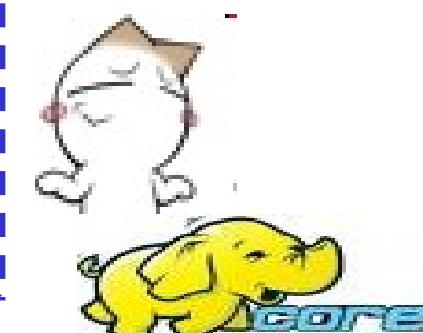
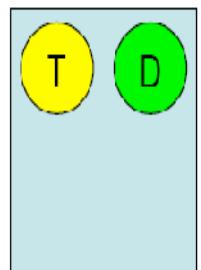
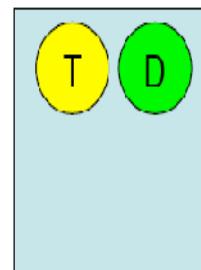
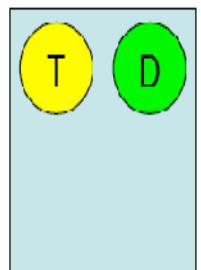
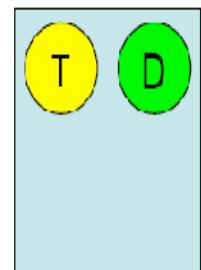
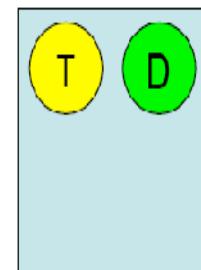
HTTP Monitoring UI

http://x.x.x.1:50070



**conf/slaves:**

Node3  
.....  
NodeN

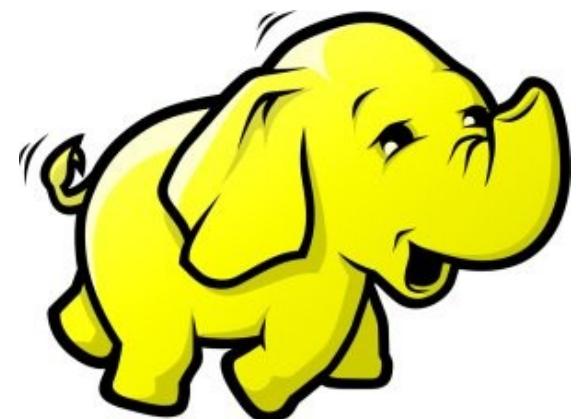




# Hadoop 簿集佈署工具

Hadoop Deployment Tool : SmartFog and DRBL

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Programmer v.s. System Admin.



Source:<http://www.funnyjunksite.com/wp-content/uploads/2007/08/programmer.jpg>



Source:  
<http://www.sysadminday.com/images/people/136-3697.JPG>



**PART 1 :**

# PC Cluster 101

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



Powered by **DRBL**



At First, We have  $4 + 1$  PC Cluster

It'd better be  
 $2^n$



Manage  
Scheduler

**Then, We connect 5 PCs with  
Gigabit Ethernet Switch**



**GiE Switch**



**10/100/1000  
MBps**



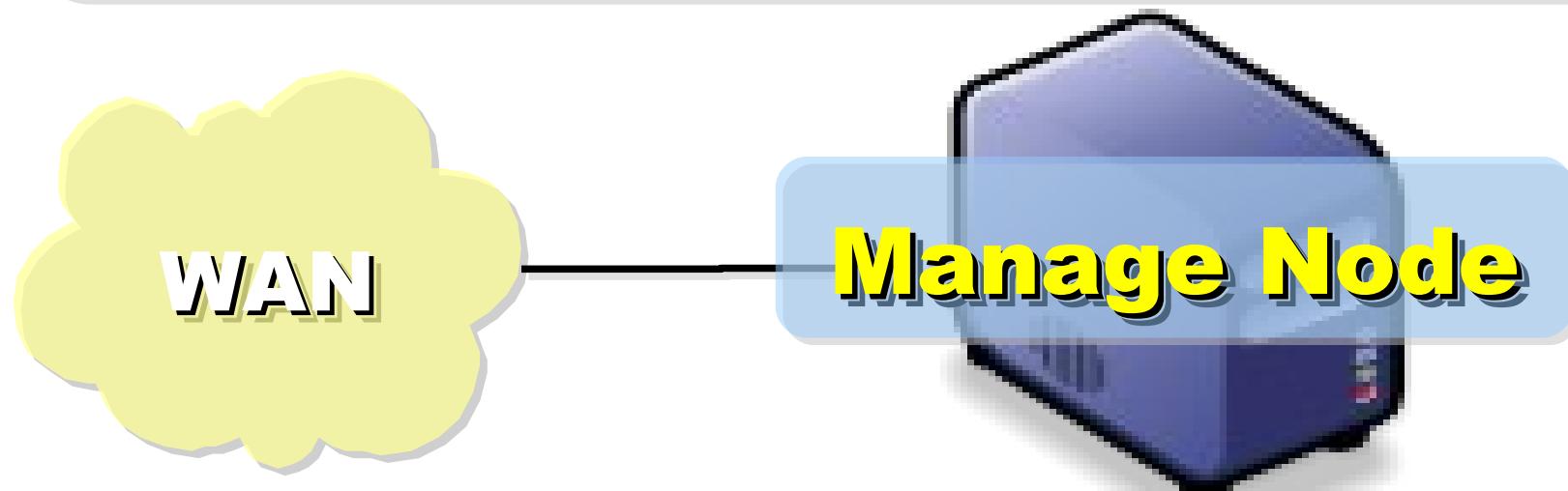
**WAN**

**Add 1 NIC  
for WAN**



## **Compute Nodes**

**4 Compute Nodes will communicate via LAN Switch. Only Manage Node have Internet Access for Security!**



# Basic System Setup for Cluster

## Compute Nodes

Messaging

**MPICH**

**GCC**

**Bash**

**Perl**

Account Mgmt.

**SSHD**

**NIS**

**YP**

**GNU Libc**



**Kernel Module**

**Linux Kernel**

**Boot Loader**

**On Manage Node,  
We need to install **Scheduler** and  
**Network File System** for sharing  
Files with Compute Node**

Job Mgmt.

**OpenPBS**

File Sharing

**NFS**



Messaging

**MPICH**

**GCC**

**Bash**

**Perl**

Account Mgmt.

**SSHD**

**NIS**

**YP**

**GNU Libc**



# Challenges of Cluster Computing

- **Hardware**
  - **Ethernet Speed / PC Density**
  - **Power / Cooling / Heat**
  - **Network and Storage Architecture**
- **Software**
  - **Job Scheduler ( Cluster level )**
  - **Account Management**
  - **File Sharing / Package Management**
- **Limitation**
  - **Shared Memory**
  - **Global Memory Management**

# Common Method to deploy Cluster



**1. Setup one  
Template  
machine**

**2. Cloning  
to  
multiple  
machine**



**3. Configure  
Settings**



**4. Install  
Job  
Scheduler**



**5. Running  
Benchmark**

# **Challenges of Common Method**

**Add New User Account ?**

**Upgrade Software ?**

**How to share user data ?**

**Configuration Synchronization**

# How to deploy 4000+ Nodes ????

資料標題 : Scaling Hadoop to 4000 nodes at Yahoo!

資料日期 : September 30, 2008

Total Nodes	4000
Total cores	30000
Data	16PB

	500-node cluster		4000-node cluster	
	write	read	write	read
number of files	990	990	14,000	14,000
file size (MB)	320	320	360	360
total MB processes	316,800	316,800	5,040,000	5,040,000
tasks per node	2	2	4	4
avg. throughput (MB/s)	5.8	18	40	66

# Advanced Methods to deploy Cluster

- **SSI ( Single System Image )**
  - **Multiple PCs as Single Computing Resources**
  - **Image-based**
    - **homogeneous**
    - **ex. SystemImager, OSCAR, Kadeploy**
  - **Package-based**
    - **heterogeneous**
    - **easy update and modify packages**
    - **ex. FAI, DRBL**
- **Other deploy tools**
  - **Rocks : RPM only**
  - **cfengine : configuration engine**

# Comparison of Cluster Deploy Tools

	<b>Distribution</b>	<b>Support Diskless/Sysmless</b>	<b>Type</b>	<b>Node configuration tools</b>	<b>Cluster management tools</b>	<b>Database installation</b>
<b>System Imager</b>	<b>ALL</b>	<b>Yes</b>	<b>Image</b>	<b>Yes</b>	<b>No</b>	<b>No</b>
<b>OSCAR</b>	<b>RPM-based</b>	<b>Yes</b>	<b>Image</b>	<b>Yes</b>	<b>Yes</b>	<b>No</b>
<b>Kadeploy</b>	<b>ALL</b>	<b>No</b>	<b>Image</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>
<b>DRBL</b>	<b>ALL</b>	<b>Yes</b>	<b>Package</b>	<b>Yes</b>	<b>Yes</b>	<b>No</b>
<b>FAI</b>	<b>Debian-Based</b>	<b>Yes</b>	<b>Package</b>	<b>Yes</b>	<b>No</b>	<b>No</b>



**PART 2-1 :**

# Hadoop Deployment Tool

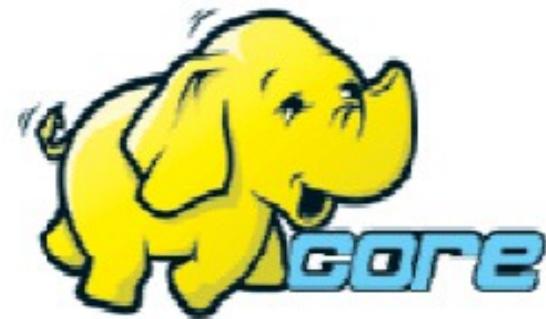
**Jazz Wang**

**Yao-Tsung Wang**

**jazz@nchc.org.tw**



Powered by **DRBL**



- Make Hadoop deployment *agile*
- Integrate with dynamic cluster deployments

Source: Deploying hadoop with smartfrog

[http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)

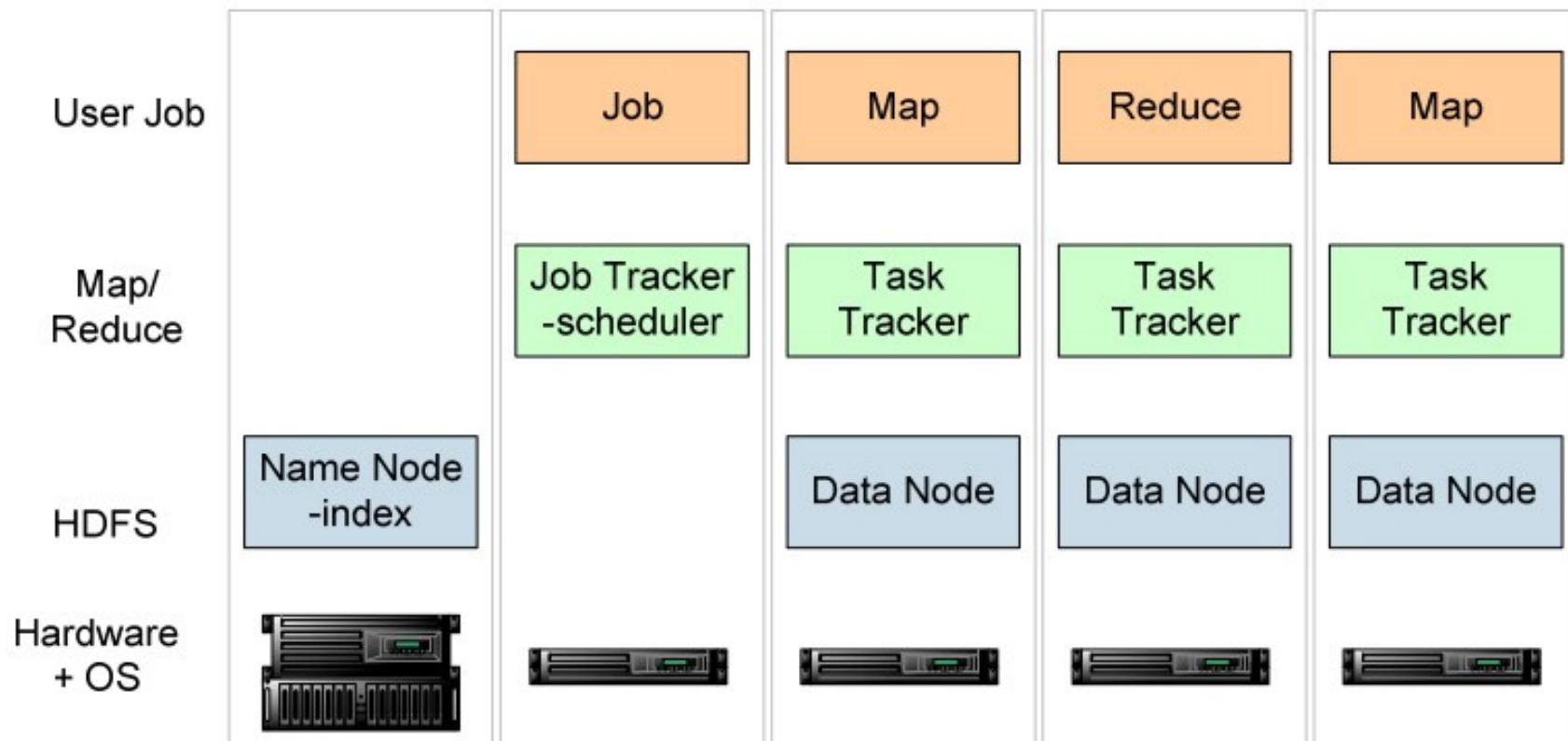
12 June 2008

# SmartFrog - HPLabs' CM tool

- Language for describing systems to deploy
  - everything from datacentres to test cases
- Runtime to create *components* from the model
- Components have a lifecycle
- LGPL Licensed, Java 5+

<http://smartfrog.org/>

# Basic problem: deploying Hadoop



*one namenode, 1+ Job Tracker, many data nodes and task trackers*

Source: Deploying hadoop with smartfrog

12 [http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)



# The hand-managed cluster

- Manual install onto machines
- SCP/FTP in Hadoop zip
- copy out hadoop-site.xml and other files
- edit /etc/hosts, /etc/rc5.d, SSH keys ...
- Installation scales  $O(N)$
- Maintenance, debugging scales worse

Source: Deploying hadoop with smartfrog

12 [http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)



# The locked-down cluster

- PXE Preboot of OS images
- RedHat Kickstart to serve up (see instalinux.com)
- Maybe: LDAP to manage state, or custom RPMs

Requires:

uniform images, central LDAP service, good ops team, stable configurations, home-rolled RPMs

Source: Deploying hadoop with smartfrog

12 [http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)



# CM-tool managed cluster

## Configuration Management tools

- State Driven: observe system state, push it back into the desired state
- Workflow: apply a sequence of operations to change a machine's state
- Centralized: central DB in charge
- Decentralized: machines look after themselves

CM tools are the only way to manage big clusters

Source: Deploying hadoop with smartfrog

[http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)



# Model the system in the SmartFrog language

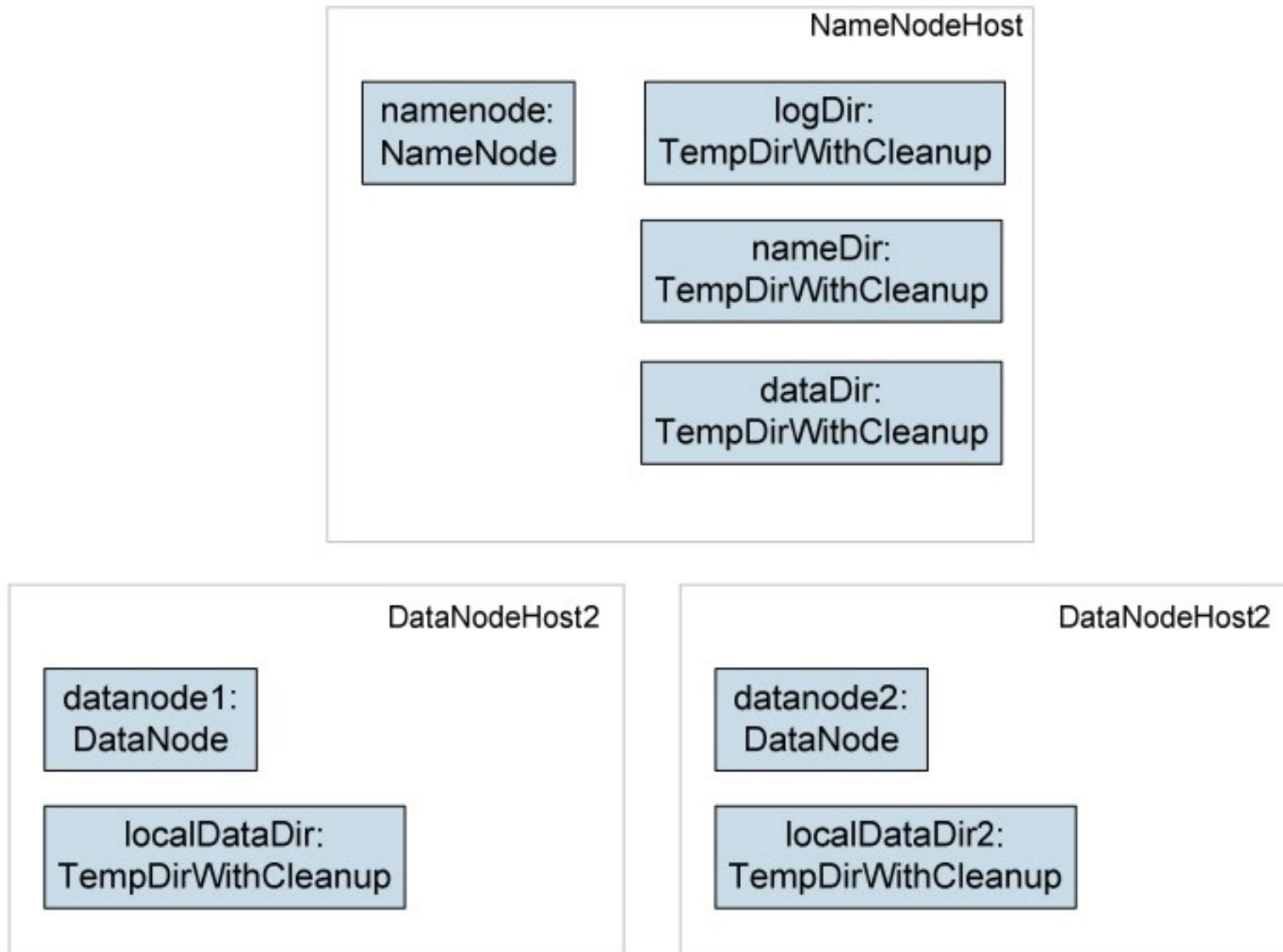
```
TwoNodeHDFS extends OneNodeHDFS {  
  
    localDataDir2 extends TempDirwithCleanup {  
  
    }  
  
    datanode2 extends datanode {  
        dataDirectories [LAZY localDataDir2];  
        dfs.datanode.https.address "https://localhost:0";  
    }  
}
```

Inheritance, cross-referencing, templating

Source: Deploying hadoop with smartfrog  
12 [http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)



# The runtime deploys the model



Source: Deploying hadoop with smartfrog

[http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)

# Steps to deployability

1. Configure Hadoop from an SmartFrog description
2. Write components for the Hadoop nodes
3. Write the functional tests
4. Add *workflow* components to work with the filesystem; submit jobs
5. Get the tests to pass

Source: Deploying hadoop with smartfrog

12 [http://people.apache.org/~stevel/slides/deploying\\_hadoop\\_with\\_smartfrog.pdf](http://people.apache.org/~stevel/slides/deploying_hadoop_with_smartfrog.pdf)





**PART 2-2 :**

## Introduction to DRBL

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



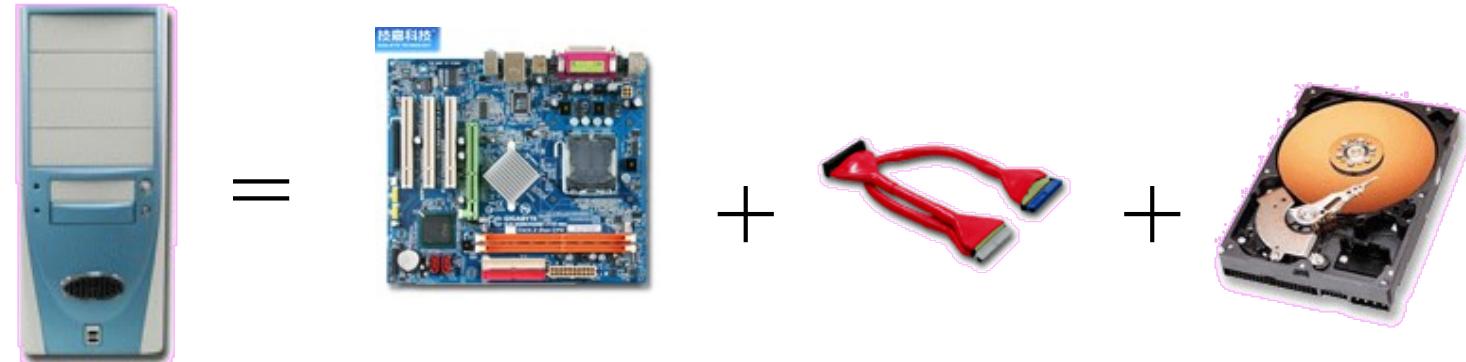
Powered by **DRBL**

# What is DRBL ??

- **Diskless Remote Boot in Linux**
- Network is cheap, and our time is expansive
- In simple words, DRBL is .....
  - Replace IDE/SATA cable with network cable
  - 40+ student PCs connected to one DRBL server



**Diskfull  
PC**



**Diskless  
PC**



**Server**

**1st, We install Base System of  
GNU/Linux on Management Node.**

**You can choose:**

**Redhat, Fedora, CentOS, Mandriva,  
Ubuntu, Debian, ...**



**2nd, We Install DRBL package and configure it as DRBL Server.**  
There are lots of service needed:  
**SSHD, DHCPD, TFTPD, NFS Server,**  
**NIS Server, YP Server ...**

Network Booting

Account Mgmt.

NFS

TFTPD

DHCPD

SSHD

NIS

YP

Perl

Bash

GNU Libc

**DRBL Server**

based on existing  
Open Source and  
keep Hacking!



Kernel Module

Linux Kernel

Boot Loader

After running “**drblsrv -j**” & “**drblpush -j**”, there will be **pxelinux**, **vmlinuz-pxe**, **initrd-pxe** in **TFTPROOT**, and different **configuration files** for each Compute Node in **NFSROOT**

**NFS**

**TFTP****D**

**DHCP****D**

**SSH****D**

**NIS**

**YP**

**Config. Files**  
Ex. **hostname**

**initrd-pxe**

**vmlinuz-pxe**

**pxelinux**

**GNU Libc**



**Kernel Module**

**Linux Kernel**

**Boot Loader**

3nd, We enable **PXE** function in **BIOS configuration.**

**BIOS PXE**

**BIOS PXE**

**BIOS PXE**

**BIOS PXE**

**NFS**

**TFTPD**

**DHCPD**

**SSHD**

**NIS**

**YP**

**Config. Files**  
**Ex. hostname**

**GNU Libc**



**initrd-pxe**

**Kernel Module**

**vmlinuz-pxe**

**Linux Kernel**

**pxelinux**

**Boot Loader**

**While Booting, PXE will query IP address from DHCPD.**

**BIOS PXE**

**BIOS PXE**

**BIOS PXE**

**BIOS PXE**

**NFS**

**TFTP**

**DHCPD**

**SSHD**

**NIS**

**YP**

**Config. Files**  
**Ex. hostname**

**GNU Libc**



**initrd-pxe**

**Kernel Module**

**vmlinuz-pxe**

**Linux Kernel**

**pxelinux**

**Boot Loader**

**While Booting, PXE will query IP address from DHCPD.**

**IP 1**

**IP 2**

**IP 3**

**IP 4**

**NFS**

**TFTPD**

**DHCPD**

**SSHD**

**NIS**

**YP**

**Config. Files**  
**Ex. hostname**

**GNU Libc**



**initrd-pxe**

**Kernel Module**

**vmlinuz-pxe**

**Linux Kernel**

**pxelinux**

**Boot Loader**

**After PXE get its IP address, it will download booting files from TFTPD.**

**IP 1**

**IP 2**

**IP 3**

**IP 4**

**NFS**

**TFTPD**

**DHCPD**

**SSHD**

**NIS**

**YP**

**Config. Files**  
**Ex. hostname**

**initrd-pxe**

**vmlinuz-pxe**

**pxelinux**

**GNU Libc**



**Kernel Module**

**Linux Kernel**

**Boot Loader**



**Config. Files**  
Ex. hostname

**initrd-pxe**

**vmlinuz-pxe**

**pxelinux**

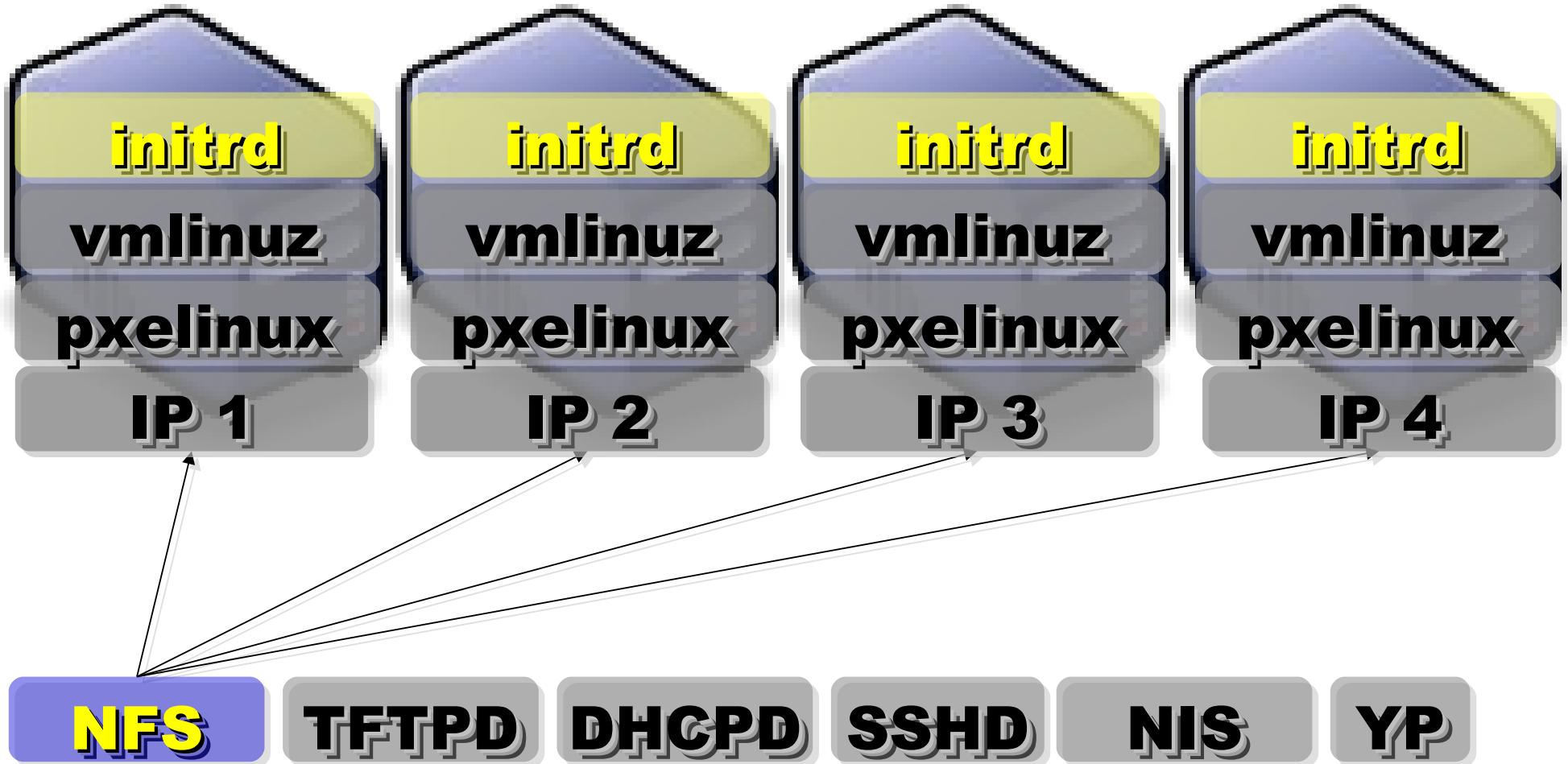
**GNU Libc**



**Kernel Module**

**Linux Kernel**

**Boot Loader**



**After downloading booting files,  
scripts in `initrd-pxe` will config  
`NFSROOT` for each Compute Node.**

## **Config. 1**

**initrd**

**vmlinuz**

**pxelinux**

**IP 1**

## **Config. 2**

**initrd**

**vmlinuz**

**pxelinux**

**IP 2**

## **Config. 3**

**initrd**

**vmlinuz**

**pxelinux**

**IP 3**

## **Config. 4**

**initrd**

**vmlinuz**

**pxelinux**

**IP 4**

**NFS**

**TFTPD**

**DHCPD**

**SSHD**

**NIS**

**YP**

**Config. Files**  
**Ex. hostname**

**initrd-pxe**

**vmlinuz-pxe**

**pxelinux**

**GNU Libc**



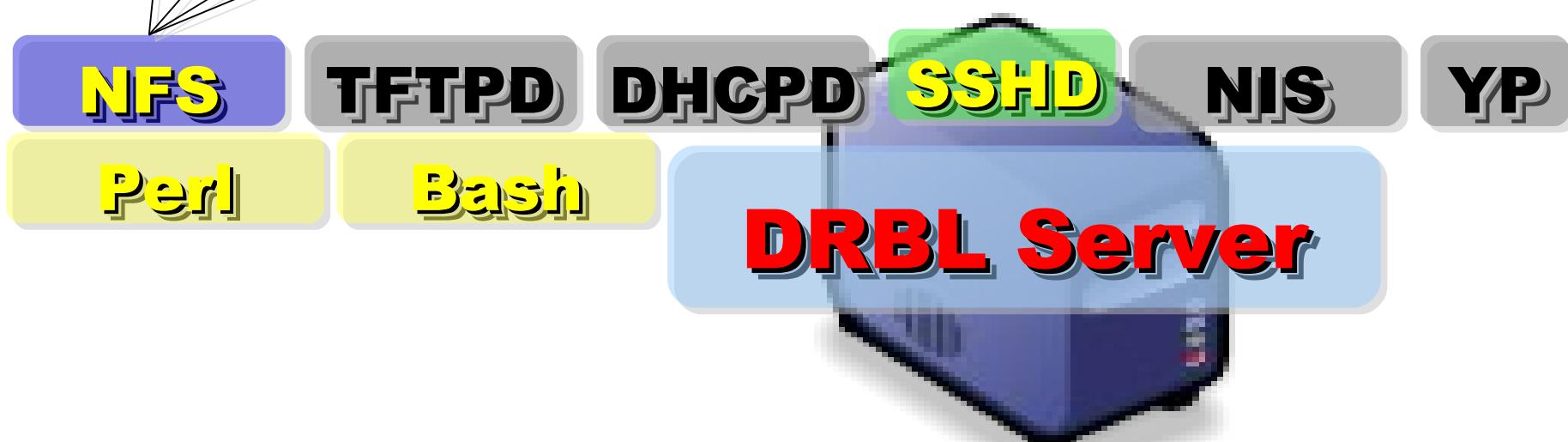
**Kernel Module**

**Linux Kernel**

**Boot Loader**



**Applications and Services will also  
deployed to each Compute Node  
via NFS ....**





**SSHD**

**SSHD**

**SSHD**

**SSHD**

**With the help of NIS and YP,  
You can login each Compute Node  
with the Same ID / PASSWORD  
stored in DRBL Server!**

**SSH Client**

**NFS**

**TFTPD**

**DHCPD**

**SSHD**

**NIS**

**YP**

**DRBL Server**



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



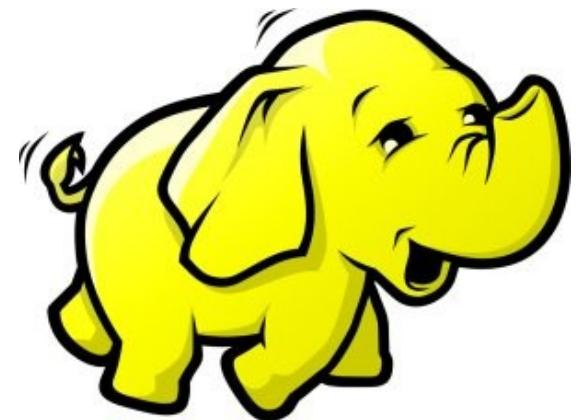
Powered by **DRBL**



# Hadoop 相關計畫

## Hadoop Ecosystem

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**





**Hadoop 只支援用 Java 開發嘛？**  
**Is Hadoop only support Java ?**

總不能全部都重新設計吧？如何與舊系統相容？

**Can Hadoop work with existing software ?**



可以跟資料庫結合嘛？

**Can Hadoop work with Databases ?**



開發者們有聽到大家的需求.....

**Yes, we hear the feedback of developers ...**



# Is Hadoop only support Java ?

- Although the Hadoop framework is implemented in Java™, **Map/Reduce applications need not be written in Java.**
- **Hadoop Streaming** is a utility which allows users to create and run jobs with any **executables** (e.g. shell utilities) as the mapper and/or the reducer.
- **Hadoop Pipes** is a SWIG-compatible C++ API to implement Map/Reduce applications (non JNI™ based).

# Hadoop Pipes (C++, Python)

- Hadoop Pipes allows **C++** code to use Hadoop DFS and map/reduce.
- The C++ interface is "swigable" so that interfaces can be generated for **python** and other scripting languages.
- For more detail, check the API Document of **org.apache.hadoop.mapred.pipes**
- You can also find example code at **hadoop-\*/src/examples/pipes**
- About the pipes C++ WordCount example code:  
**[http://wiki.apache.org/hadoop/C++WordCount](http://wiki.apache.org/hadoop/C%2B%2BWordCount)**

# Hadoop Streaming

- Hadoop Streaming is a utility which allows users to create and run Map-Reduce jobs **with any executables** (e.g. Unix shell utilities) as the mapper and/or the reducer.
- It's useful when you need to run **existing program** written in shell script, perl script or even PHP.
- Note: both the **mapper** and the **reducer** are **executables** that read the input from **STDIN** (line by line) and emit the output to **STDOUT**.
- For more detail, check the official document of **Hadoop Streaming**

# Running Hadoop Streaming

```
jazz@hadoop:~$ hadoop jar hadoop-streaming.jar -help  
10/08/11 00:20:00 ERROR streaming.StreamJob: Missing required option -input  
Usage: $HADOOP_HOME/bin/hadoop [--config dir] jar \  
      $HADOOP_HOME/hadoop-streaming.jar [options]
```

Options:

-input <path>	<b>DFS input file(s)</b> for the Map step
-output <path>	<b>DFS output directory</b> for the Reduce step
-mapper <cmd JavaClassName>	<b>The streaming command to run</b>
-combiner <JavaClassName>	Combiner has to be a Java class
-reducer <cmd JavaClassName>	<b>The streaming command to run</b>
-file <file>	File/dir to be shipped in the Job jar file
-dfs <h:p> local	Optional. Override DFS configuration
-jt <h:p> local	Optional. Override JobTracker configuration
-additionalconfspec specfile	Optional.
-inputformat <b>TextInputFormat(default)</b>   SequenceFileAsTextInputFormat   JavaClassName	Optional.
-outputformat <b>TextOutputFormat(default)</b>   JavaClassName	Optional.

# Hadoop Streaming with shell commands (1)

```
hadoop:~$ hadoop fs -rmr input output  
hadoop:~$ hadoop fs -put /etc/hadoop/conf input  
hadoop:~$ hadoop jar hadoop-streaming.jar -input  
input -output output -mapper /bin/cat  
-reducer /usr/bin/wc
```

# Hadoop Streaming with shell commands (2)

```
hadoop:~$ echo "sed -e \"s/ /\n/g\" | grep ." > streamingMapper.sh  
hadoop:~$ echo "uniq -c | awk '{print \$2 \"\t\" \$1}'" > streamingReducer.sh  
hadoop:~$ chmod a+x streamingMapper.sh  
hadoop:~$ chmod a+x streamingReducer.sh  
hadoop:~$ hadoop fs -put /etc/hadoop/conf input  
hadoop:~$ hadoop jar hadoop-streaming.jar -input  
input -output output -mapper streamingMapper.sh  
-reducer streamingReducer.sh -file  
streamingMapper.sh -file streamingReducer.sh
```

# There are several Hadoop subprojects

Apache > Hadoop >

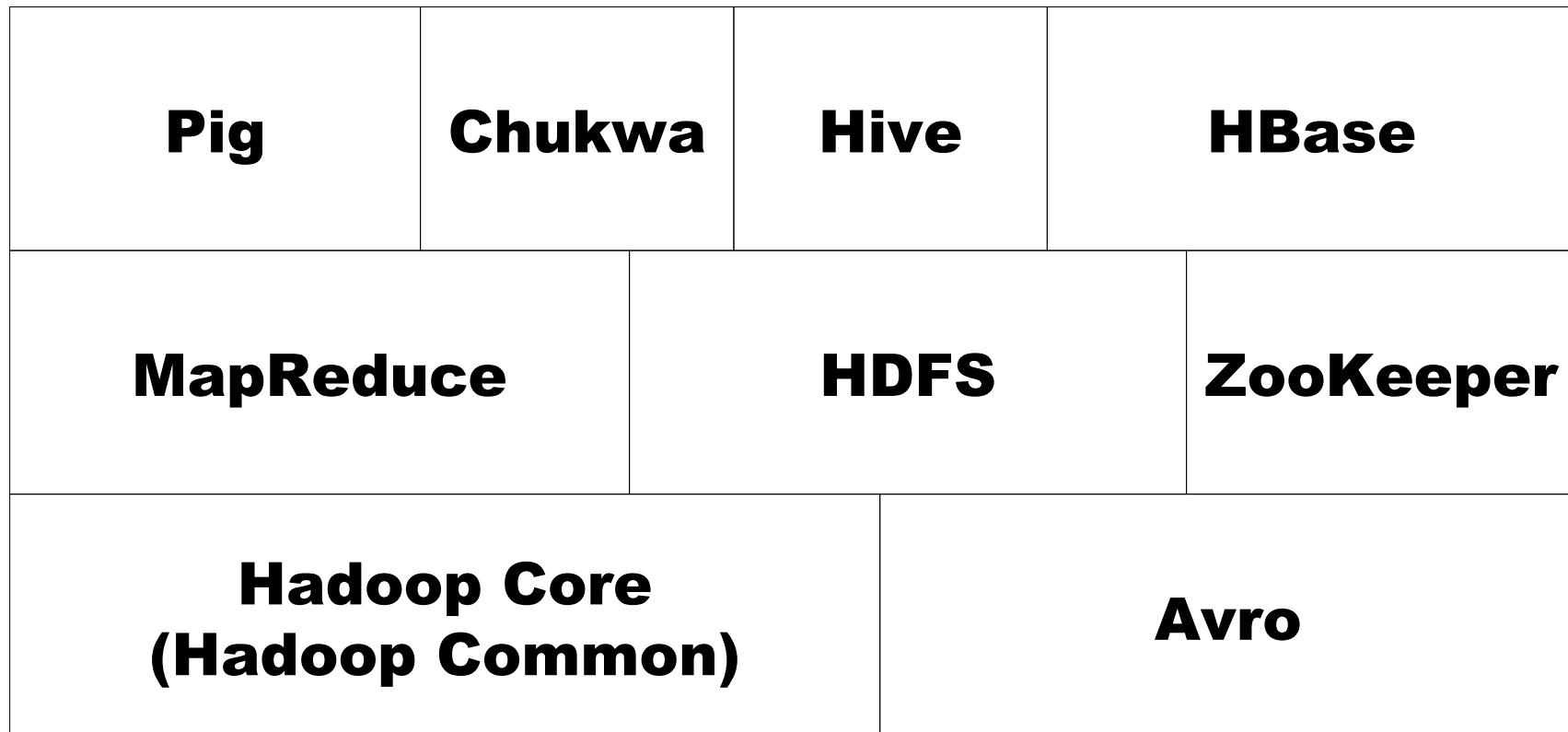


- **Hadoop Common:** The common utilities that support the other Hadoop subprojects.
- **HDFS:** A distributed file system that provides high throughput access to application data.
- **MapReduce:** A software framework for distributed processing of large data sets on compute clusters.

# Other Hadoop related projects

- **Chukwa**: A data collection system for managing large distributed systems.
- **HBase**: A scalable, distributed database that supports structured data storage for large tables.
- **Hive**: A data warehouse infrastructure that provides data summarization and ad hoc querying.
- **Pig**: A high-level data-flow language and execution framework for parallel computation.
- **ZooKeeper**: A high-performance coordination service for distributed applications.

# Hadoop Ecosystem



Source: *Hadoop: The Definitive Guide*

# Avro

- Avro is a **data serialization system**.
- It provides:
  - *Rich data structures*.
  - *A compact, fast, binary data format*.
  - *A container file, to store persistent data*.
  - *Remote procedure call (RPC)*.
  - *Simple integration with dynamic languages*.
- Code generation is not required to read or write data files nor to use or implement RPC protocols. Code generation as an optional optimization, only worth implementing for statically typed languages.
- For more detail, please check the official document:  
<http://avro.apache.org/docs/current/>



# Zoo Keeper



- <http://hadoop.apache.org/zookeeper/>
- ZooKeeper is a **centralized service** for maintaining **configuration** information, naming, providing distributed **synchronization**, and providing group services. All of these kinds of services are used in some form or another by distributed applications.
- *Each time they are implemented there is a lot of work that goes into fixing the bugs and race conditions that are inevitable. Because of the difficulty of implementing these kinds of services, applications initially usually skimp on them ,which make them brittle in the presence of change and difficult to manage. Even when done correctly, different implementations of these services lead to management complexity when the applications are deployed.*

# Pig

- <http://hadoop.apache.org/pig/>
- Pig is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs.
- Pig's infrastructure layer consists of a compiler that produces sequences of Map-Reduce programs
- Pig's language layer currently consists of a textual language called Pig Latin, which has the following key properties:
  - Ease of programming
  - Optimization opportunities
  - Extensibility



# Hive

- <http://hadoop.apache.org/hive/>
- Hive is a **data warehouse** infrastructure built on top of Hadoop that provides tools to enable easy **data summarization**, **adhoc querying** and analysis of large datasets data stored in Hadoop files.
- **Hive QL** is based on SQL and enables users familiar with SQL to query this data.



# Chukwa

- <http://hadoop.apache.org/chukwa/>
- Chukwa is an open source **data collection system** for monitoring large distributed systems.
- built on top of HDFS and Map/Reduce framework
- includes a flexible and powerful toolkit for displaying, monitoring and analyzing results to make the best use of the collected data.



# Mahout

- <http://mahout.apache.org/>
- Mahout is a scalable **machine learning libraries**.
- implemented on top of Apache Hadoop using the map/reduce paradigm.
- Mahout currently has
  - Collaborative Filtering
  - User and Item based recommenders
  - K-Means, Fuzzy K-Means clustering
  - Mean Shift clustering
  - More ...

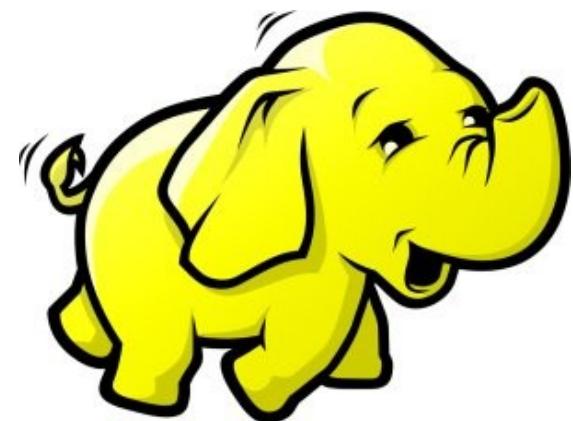




# HBase 雲端資料庫

## Introduction to HBase

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# It's all about SCALE!!

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194  
Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194  
Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194  
Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194  
Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194  
Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

訂購歷史紀錄



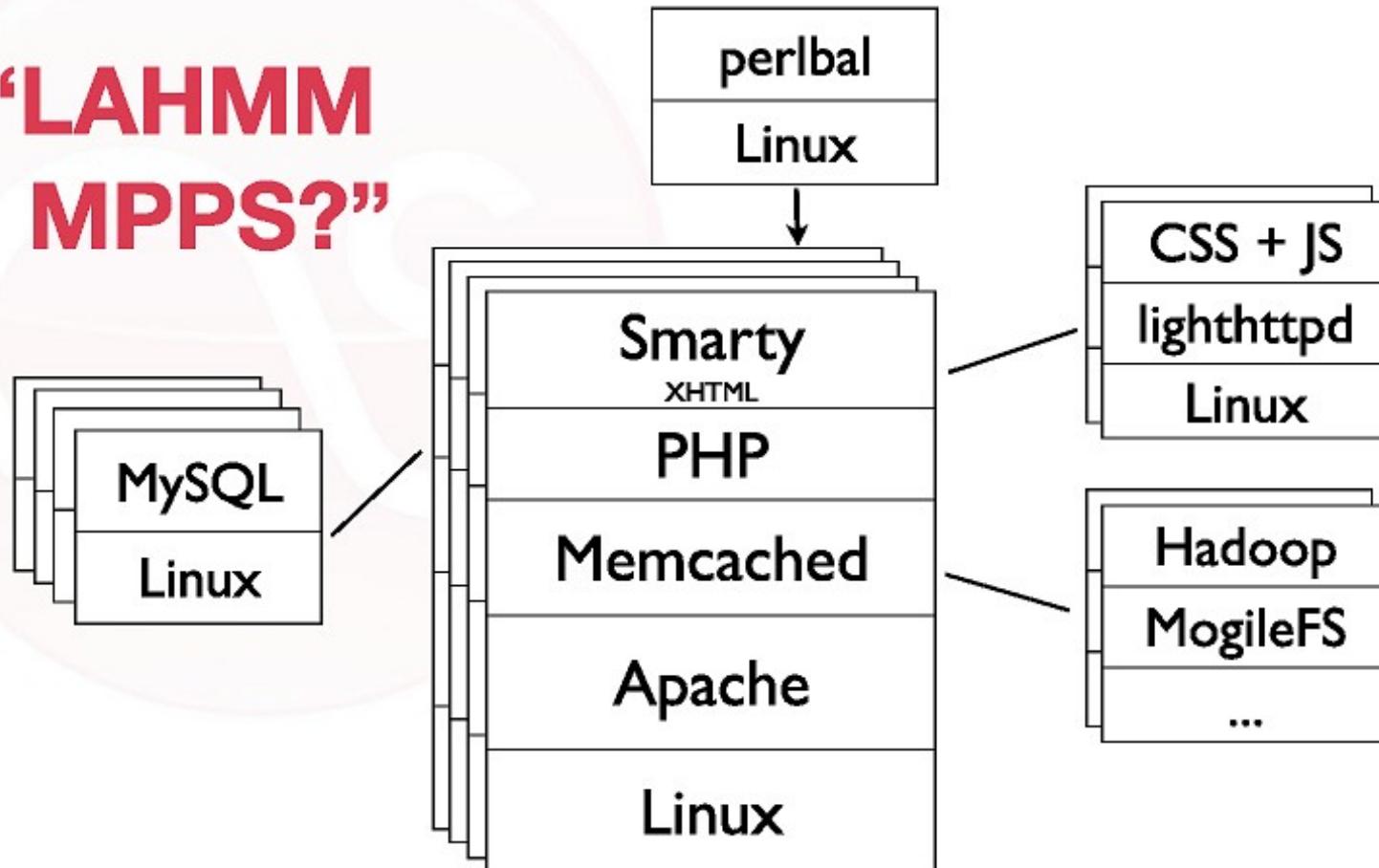
denied in /usr/local/apache/htdocs/include2007/oracle/db\_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public\_html/./logs/oracle\_error\_log.txt)

**Warning:** fopen(/home/dodgers/public\_html/./logs/oracle\_error\_log.txt) [function.fopen]: failed to open stream: Permission

# How to scale up web service in the past ?

**“LAHMM  
MPPS?”**



Where we can go: horizontal LAMP scaling example

2. A few definitions



The Social Music Revolution  
© Last.fm 2007. For internal use only.

# Tools used by large scale websites

- Perlbal - <http://www.danga.com/perlbal/>

- ◆ 多個網頁伺服器的負載平衡
- ◆ Load balancer

- MogileFS - <http://www.danga.com/mogilefs/>

- ◆ 分散式檔案系統
- ◆ Distributed File System for small files
- ◆ 有公司認為 MogileFS 比起 Hadoop 適合拿來處理小檔案

- memcached - <http://memcached.org/>

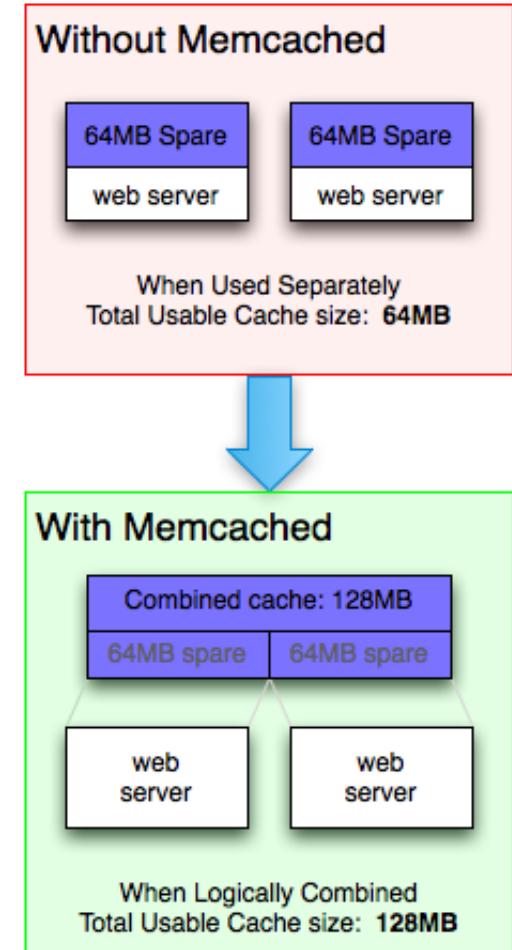
- ◆ 共享記憶體 ??
- ◆ Share Memory
- ◆ 把資料庫或經常讀取的部分，  
用記憶體快取 (Cache) 方式存放

- Moxi - <http://code.google.com/p/moxi/>

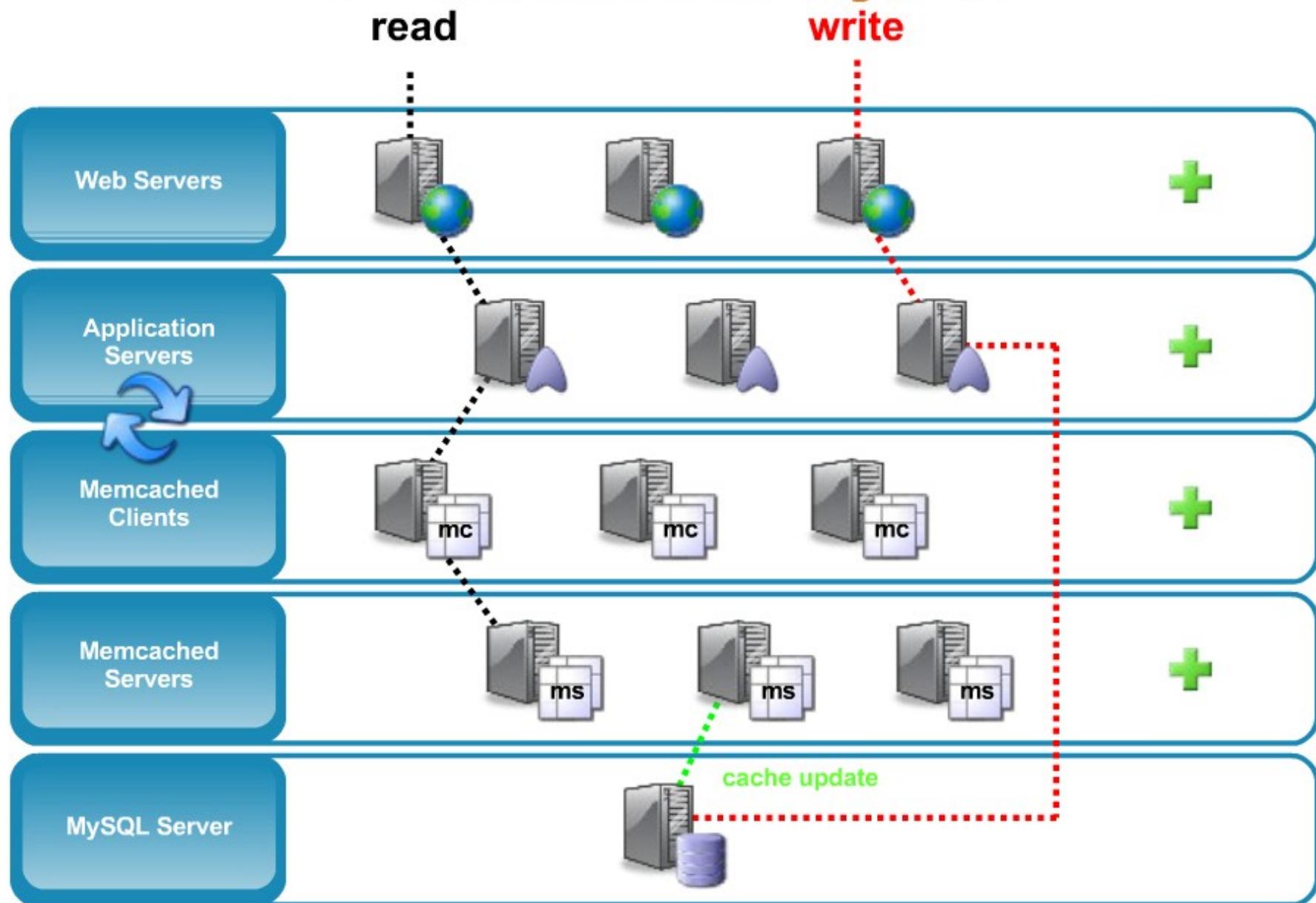
- ◆ Memcache 的 PROXY

- More Resource:

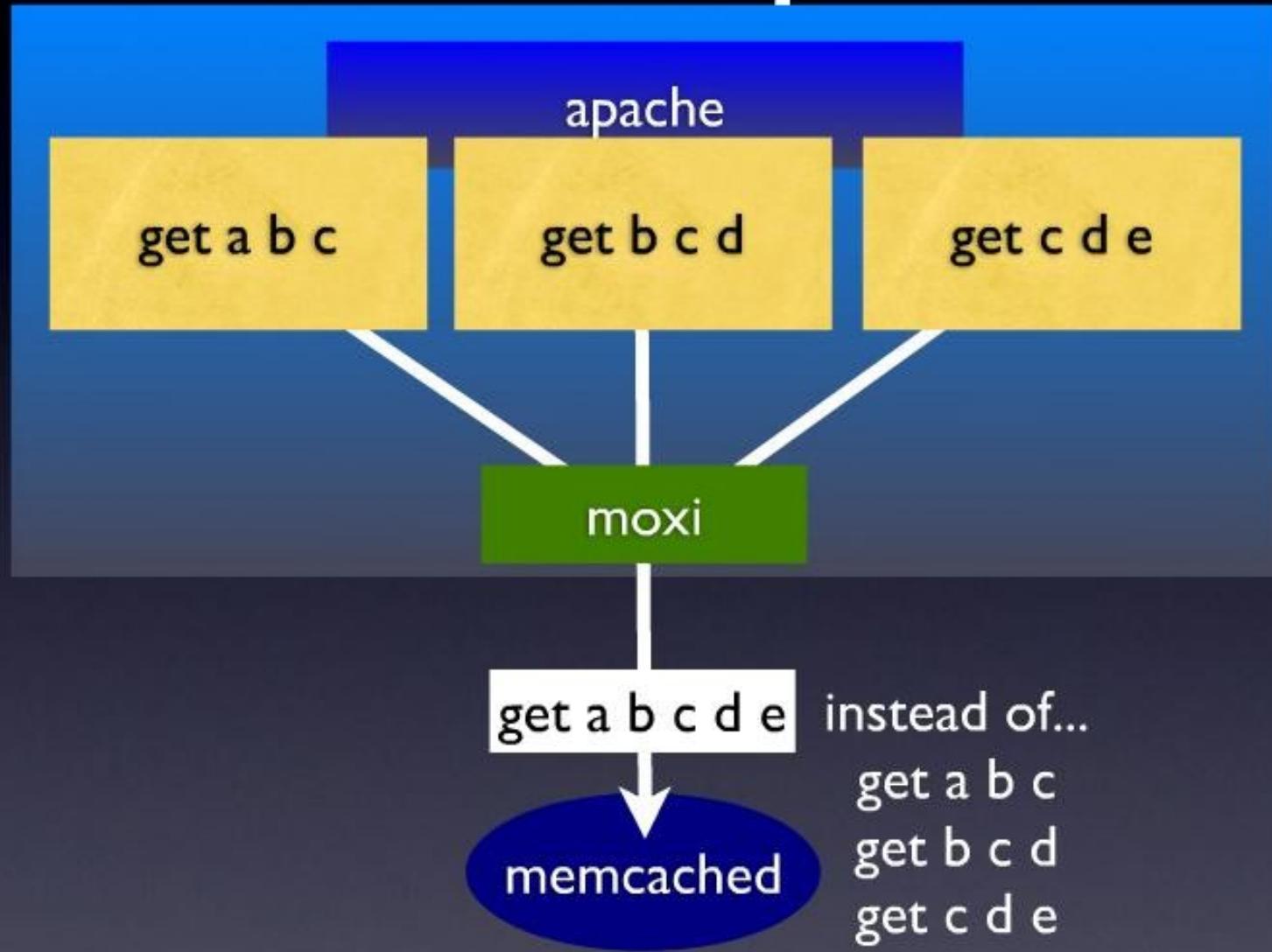
- ◆ <http://code.google.com/p/memcached/wiki/HowToLearnMoreScalability>
- ◆ <http://www.slideshare.net/techdude/scalable-web-architectures-common-patterns-and-approaches>



# Memcached & MySQL

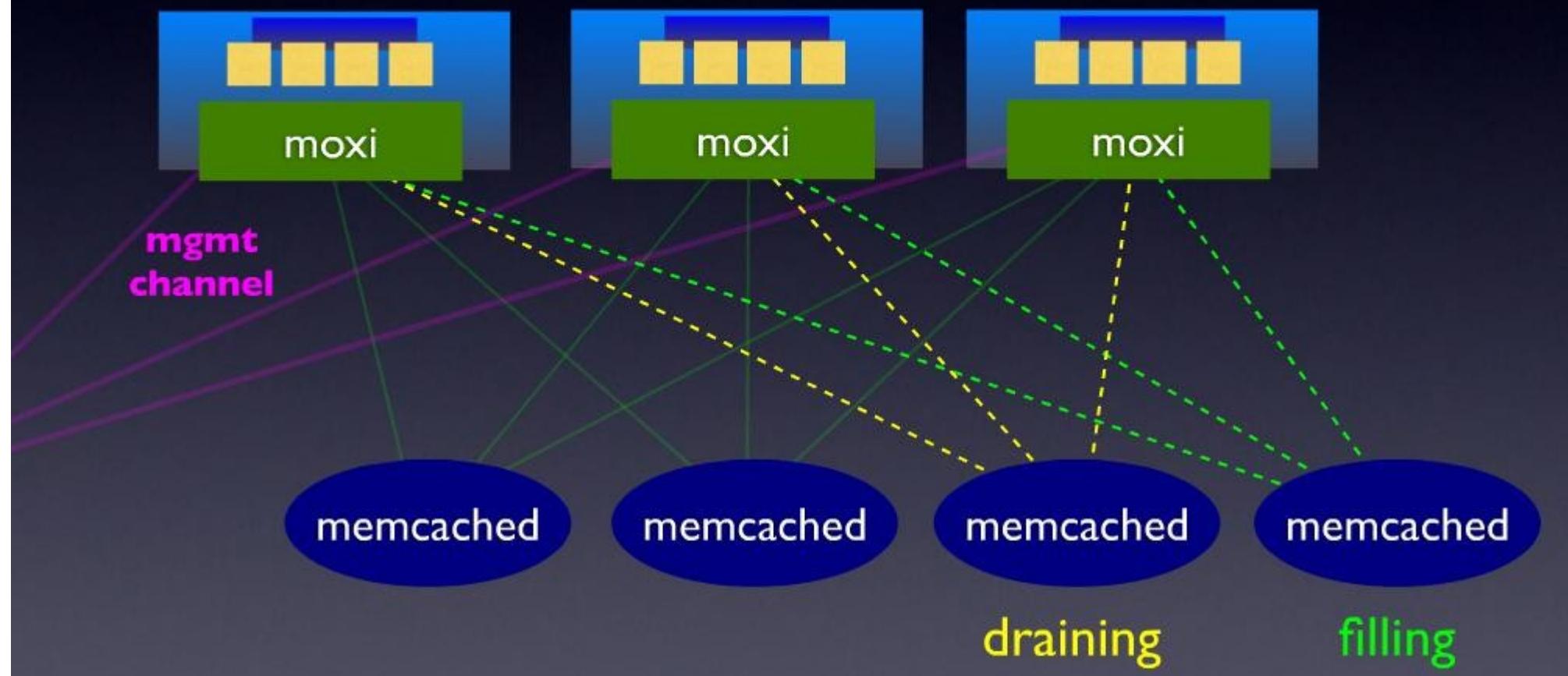


# GET de-duplication



# draining and filling

lazily migrate items from old server to new server



# HBase is ..

- HBase is a distributed **column-oriented database** built on top of HDFS.
- A distributed data store that can scale horizontally to 1,000s of commodity servers and **petabytes** of indexed storage.
- Designed to operate on top of the Hadoop distributed file system (**HDFS**) or Kosmos File System (**KFS**, aka Cloudstore) for scalability, fault tolerance, and high availability.
- Integrated into the Hadoop **map-reduce** platform and paradigm.

# Benefits

- Distributed storage
- Table-like in data structure
  - multi-dimensional map
- High scalability
- High availability
- High performance

# Who use HBase

- Adobe
  - 內部使用 (Structure data)
- Kalooga
  - 圖片搜尋引擎 <http://www.kalooga.com/>
- Meetup
  - 社群聚會網站 <http://www.meetup.com/>
- Streamy
  - Migrate from MySQL to Hbase <http://www.streamy.com/>
- Trend Micro
  - 雲端掃毒架構 <http://trendmicro.com/>
- Yahoo!
  - 儲存文件 fingerprint 避免重複 <http://www.yahoo.com/>
- More - <http://wiki.apache.org/hadoop/Hbase/PoweredBy>

# Backdrop

- Started toward by Chad Walters and Jim
- 2006.11
  - Google releases paper on **BigTable**
- 2007.2
  - Initial HBase prototype created as Hadoop contrib.
- 2007.10
  - First useable HBase
- 2008.1
  - Hadoop become Apache top-level project and HBase becomes subproject
- 2008.10~
  - HBase 0.18, 0.19 released

# HBase Is Not ...

- Tables have **one primary index**, the *row key*.
- **No join operators.**
- Scans and queries can select a subset of available columns, perhaps by using a wildcard.
- There are three types of lookups:
  - Fast lookup using row key and optional timestamp.
  - Full table scan
  - Range scan from region start to end.

# HBase Is Not ... (2)

- Limited atomicity and transaction support.
  - HBase supports **multiple batched mutations of single rows** only.
  - Data is unstructured and untyped.
- No accessed or manipulated via SQL.
  - Programmatic access via Java, REST, or **Thrift APIs**.
  - Scripting via JRuby.

# Why Bigtable?

- Performance of RDBMS system is good for transaction processing but for very large scale analytic processing, the solutions are commercial, expensive, and specialized.
- Very large scale analytic processing
  - Big queries – typically range or table scans.
  - **Big databases (100s of TB)**

# Why Bigtable? (2)

- Map reduce on Bigtable with optionally Cascading on top to support some relational algebras may be a cost effective solution.
- Sharding is not a solution to scale open source RDBMS platforms
  - Application specific
  - Labor intensive (**re**)partitioning

# Why HBase ?

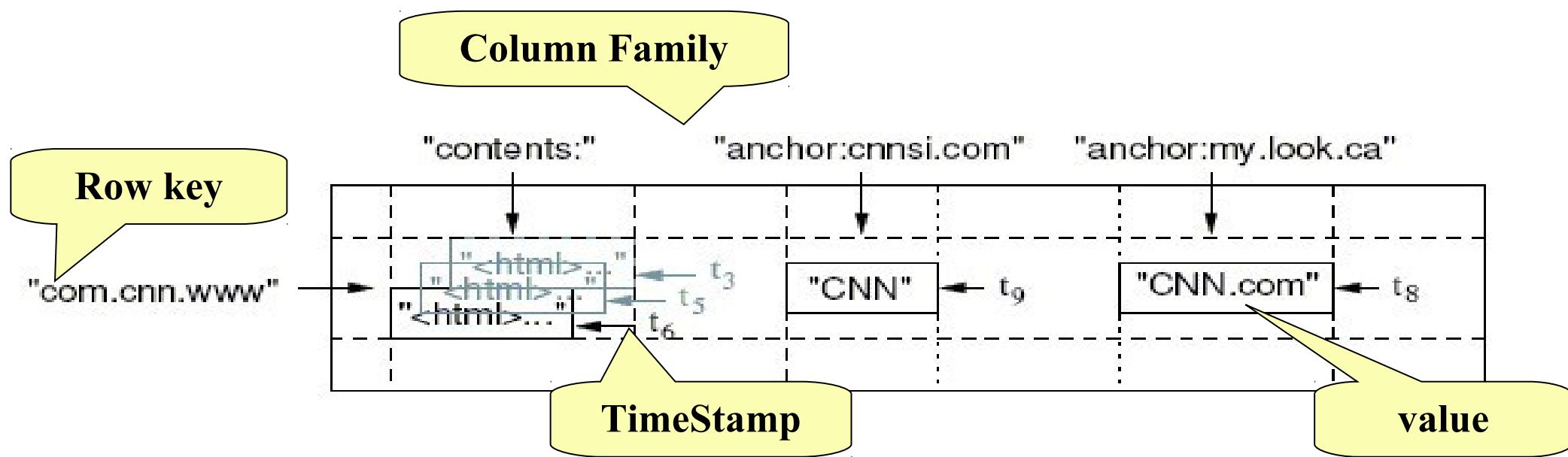
- HBase is a Bigtable clone.
- It is open source
- It has a good community and promise for the future
- It is developed on top of and has good integration for the Hadoop platform, if you are using Hadoop already.
- It has a Cascading connector.

# HBase benefits than RDBMS

- *No real indexes*
- *Automatic partitioning*
- *Scale linearly and automatically with new nodes*
- *Commodity hardware*
- *Fault tolerance*
- *Batch processing*

# Data Model

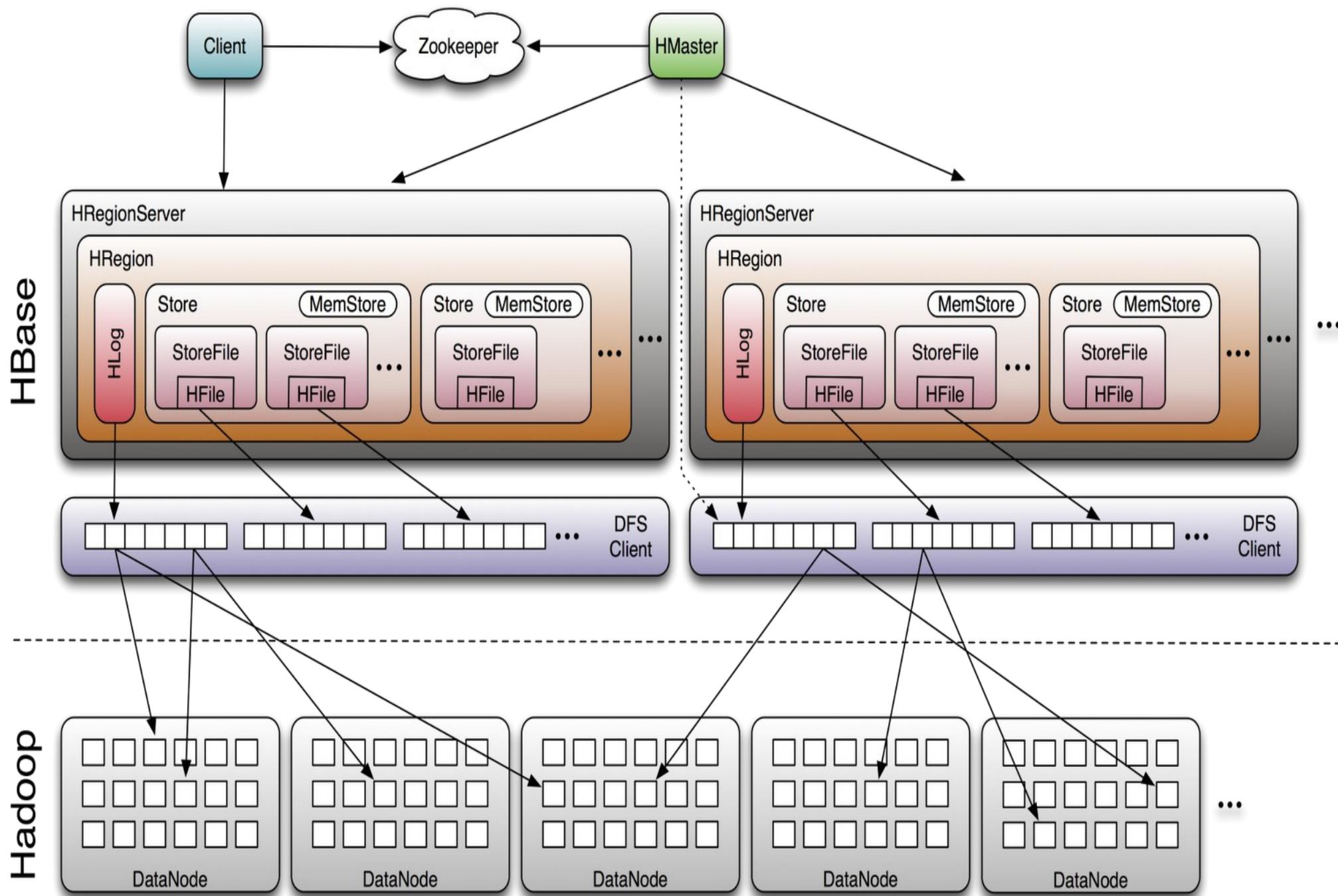
- Tables are sorted by **Row**
- Table schema only define it's *column families* .
  - Each family consists of any number of columns
  - Each column consists of any number of versions
  - Columns only exist when inserted, NULLs are free.
  - Columns within a family are sorted and stored together
- Everything except table names are byte[]
- **(Row, Family: Column, Timestamp) → Value**



# Members

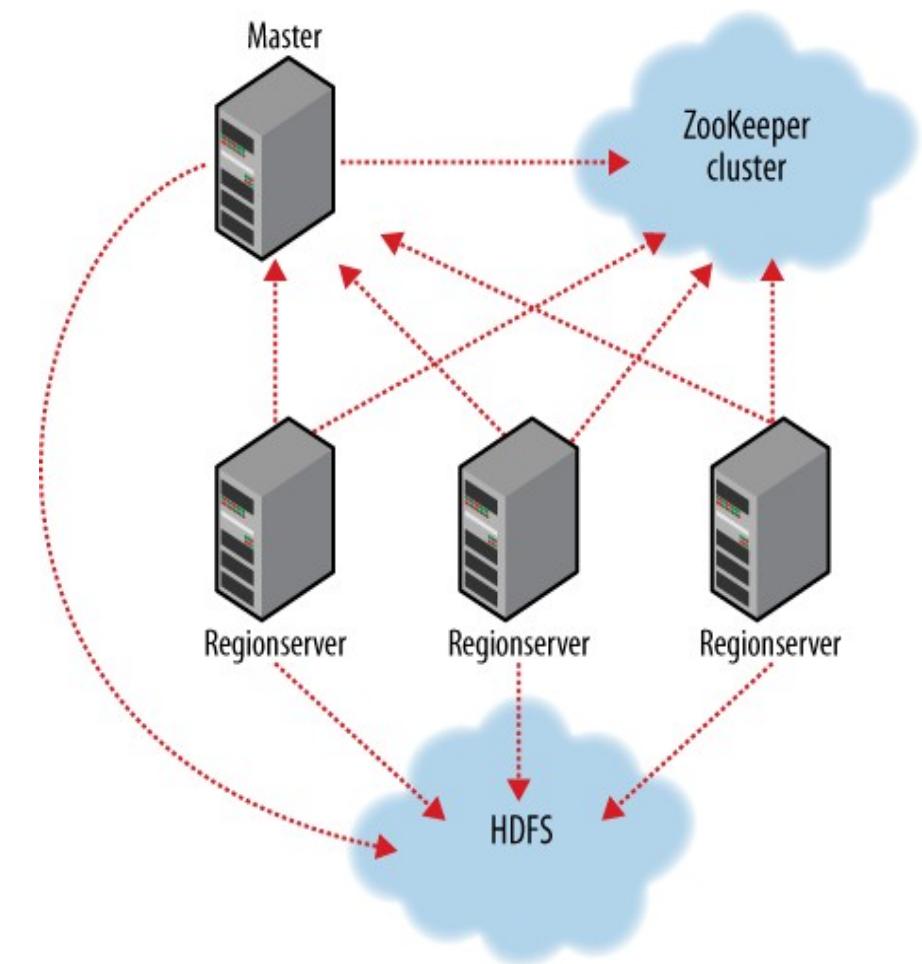
- *Master*
  - Responsible for monitoring region servers
  - Load balancing for regions
  - Redirect client to correct region servers
  - The current SPOF
- *regionserver slaves*
  - Serving requests(Write/Read/Scan) of Client
  - Send HeartBeat to Master
  - Throughput and Region numbers are scalable by region servers

# Architecture



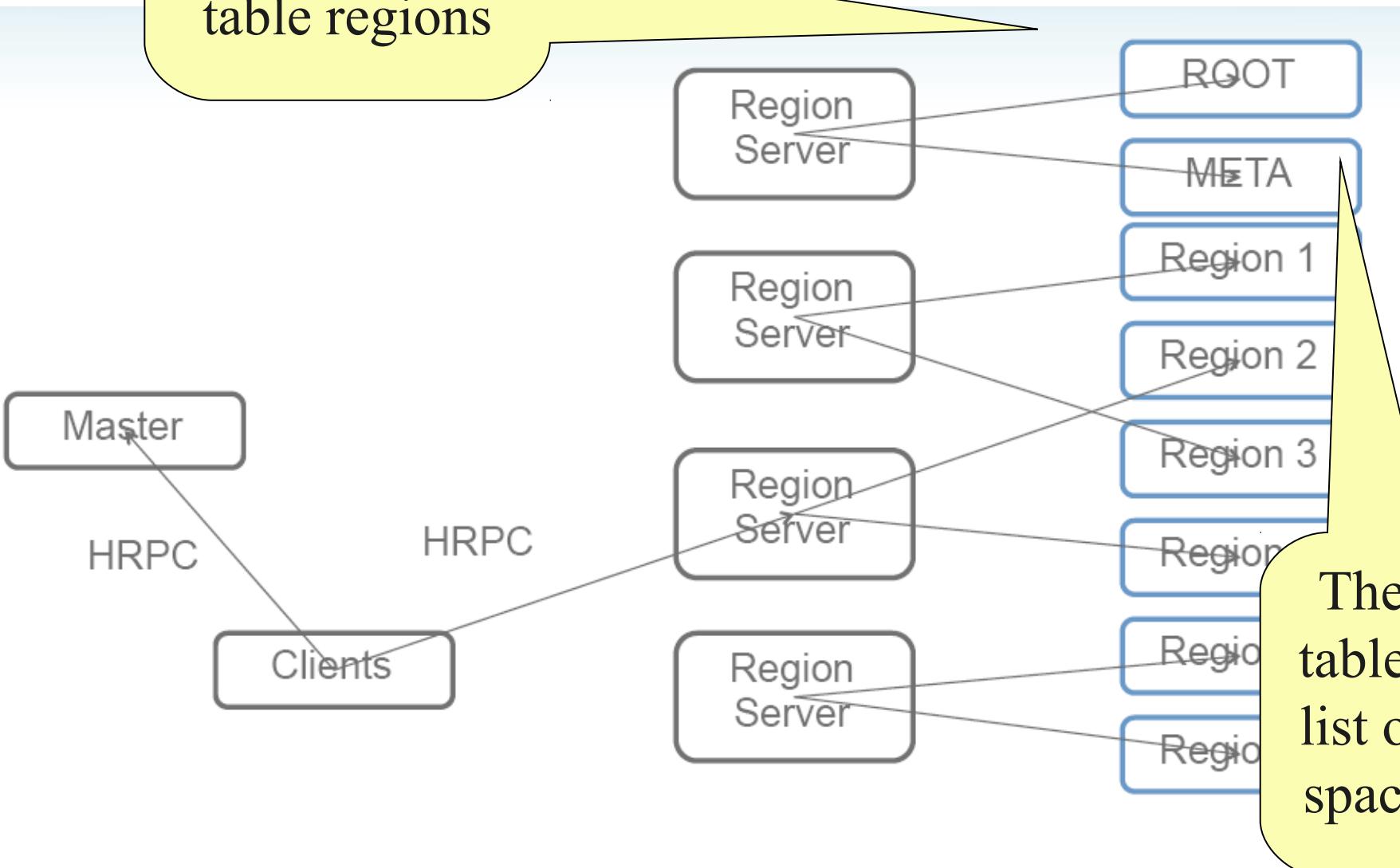
# ZooKeeper

- HBase depends on ZooKeeper (Chapter 13) and by default it manages a ZooKeeper instance as the authority on cluster state



# Operation

The -ROOT-table holds the list of .META. table regions



The .META. table holds the list of all user-space regions.



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



Powered by **DRBL**

# Introduction to Pig programming



**Yahoo Search Engineering**  
陳奕璋 ( Yiwei Chen )



# 任務！

[殺很大、插很大\(+瑤瑤寫真性感精選54P\) @ osaki's Blog :: Xuite日誌](#)

殺不用錢～殺online瑤瑤性感變裝照+精選性感寫真童顏巨乳的娃娃音美少女瑤瑤 本名：郭書瑤  
暱稱：瑤瑤 身高：155cm 體重：42kg 三圍：33E/23/33 生日：1990/7/18 ...

[blog.xuite.net/osaki99/blog/21865265 - 頁庫存檔 - 類似內容](#)

[電玩美少女瑤瑤精選影音\(ヤオヤオ童顔Fカップ爆乳美少女映画videos ...\)](#)

2008年8月31日 ... 18歲電玩少女瑤瑤半工半讀扛家計 (內有瑤瑤男友) <http://blog.xuite.net/kaiger/daily/23136438> ... 20080913 我猜嚟嚟美少女第二段2號Kiki 3號瑤瑤 ...

[blog.xuite.net/kaiger/daily/19128818 - 頁庫存檔 - 類似內容](#)

 [顯示更多來自 blog.xuite.net 的結果](#)

[jays1943 分享正妹NO.24 無名瑤瑤- 樂多日誌](#)

瑤瑤也沒有哪裡得罪你們押你們為審麼這樣罵他說害女生生氣我看你們長的很醜吧不要自以為是  
喔死網友還罵人〈死勒你要不要臉瑤瑤可是我的偶像你們最好是向一點 ...

[blog.roodo.com/jays1943/archives/6850053.html - 頁庫存檔 - 類似內容](#)



# 任務！

## 瑤瑤航空 - Powered by Discuz!

瑤瑤航空 - Discuz! Board ... 歡迎VIP旅客-魏如昀加入瑤瑤航空(2008-4-7) 歡迎VIP旅客-賴銘偉加入瑤瑤航空(2008-3-13) ... 瑶瑤家族 瑶瑤在雅虎的第一家族 瑶瑤天空部落格 林佩瑤在天空的部落格 林佩瑤 無名網誌 瑶瑤的新照片都在無名啦！無不癡齋 ...

[www.yaoyaofly.com](http://www.yaoyaofly.com) - 庫存頁面 - 更多此站結果

## 瑤瑤喵小屋~ - 無名小站

瑤瑤喵小屋~- 無名小站 Blog Album... 最近好煩煩煩，我覺得我的腦容量變小了... 好多事情消化不良 好多念頭讓我無法抉擇 (More.) goukigouki at 無名小站 at 02:39 PM post | Reply(27) | Trackback(0) | prosecute ...

[www.wretch.cc/blog/goukigouki](http://www.wretch.cc/blog/goukigouki) - 74k - 庫存頁面 - 更多此站結果

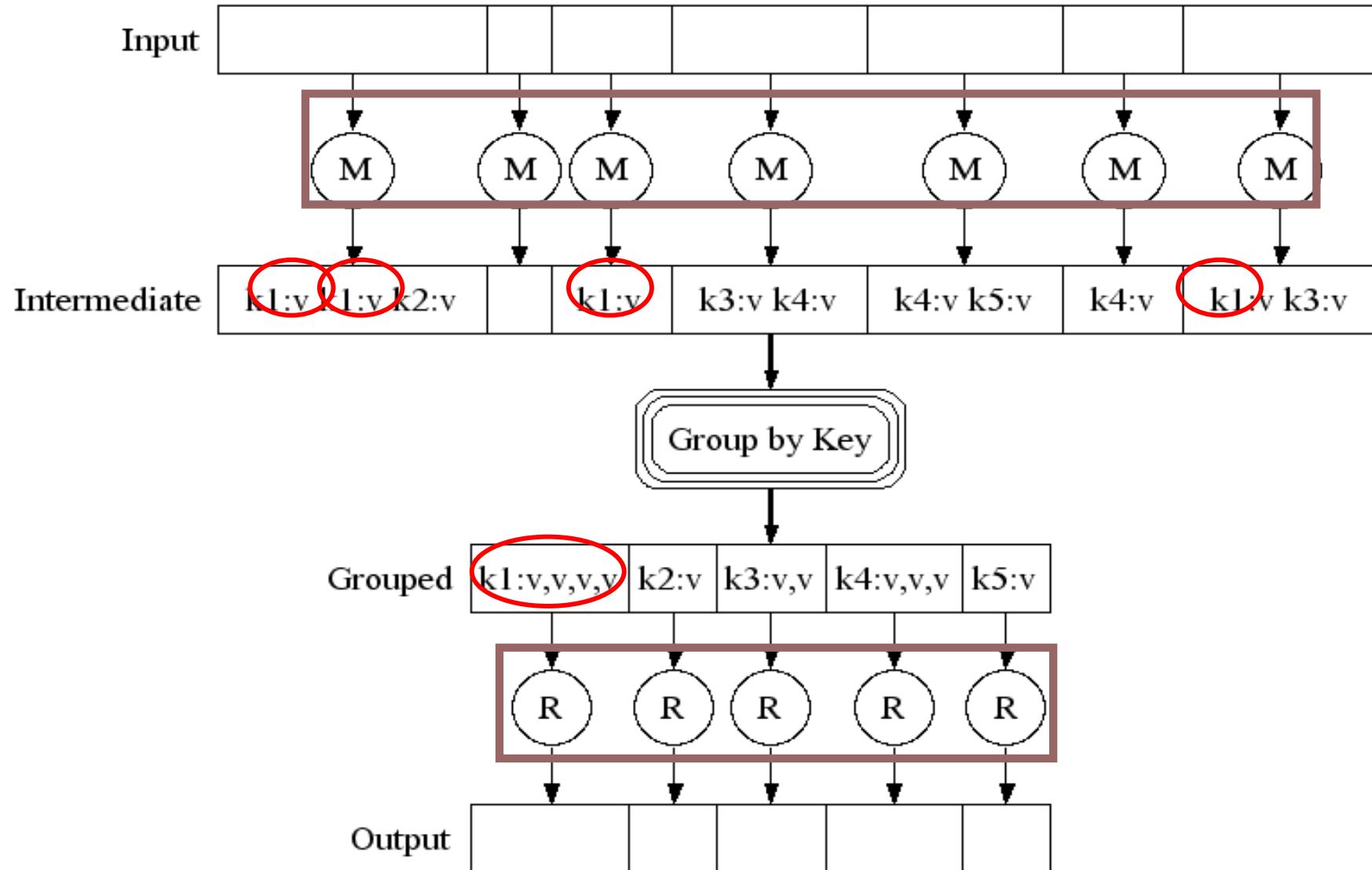


# 任務！

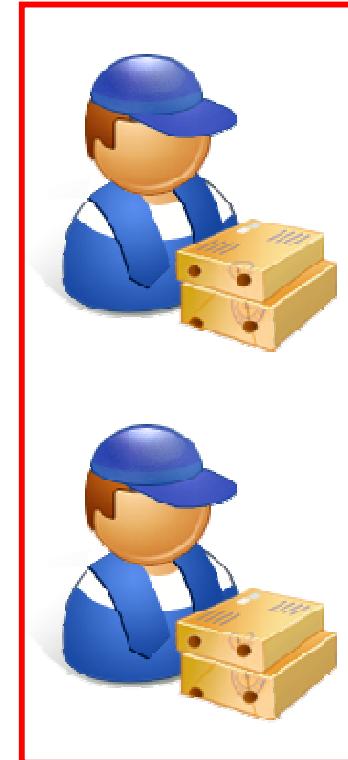
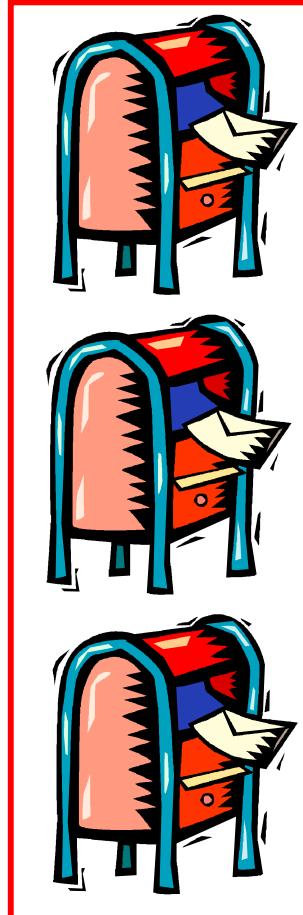
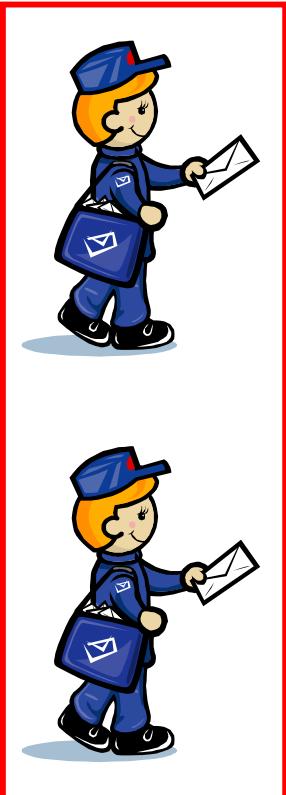
- 你怎麼知道我們放的網頁比較好？
- 你怎麼知道第一筆結果應該要多熱門？



# Hadoop Programming – Map/Reduce



# Map / Reduce



100:  
37



220:  
28

**mappers**

**reducers**

# Map-Reduce

- 全新想法
- 須分別撰寫 mappers & reducers
- 會有超級無敵霹靂多的 mapper/reducer  
要維護！

# We usually do ...

- 大部份時候：
  - filtering, projecting
  - grouping, aggregation, joining
- 今天有多少人搜尋「美國生」

# Pig (Latin)

- Procedural dataflow language (Pig Latin) for Map-Reduce
  - 很像 SQL
    - group, join, filter, sort ...
  - 人人都會 SQL

# Pig Script Example

- Top sites visited by users aged 18 to 25

```
Users = LOAD 'users.in' AS (name, age);
Fltrd = FILTER Users by age >= 18 and age <= 25;

Pages = LOAD 'pages.in' AS (user, url);

Jnd   = JOIN Fltrd BY name, Pages BY user;
Grpd  = GROUP Jnd by url;
Smmd  = FOREACH Grpd GENERATE group, COUNT(Jnd) AS
        clicks;

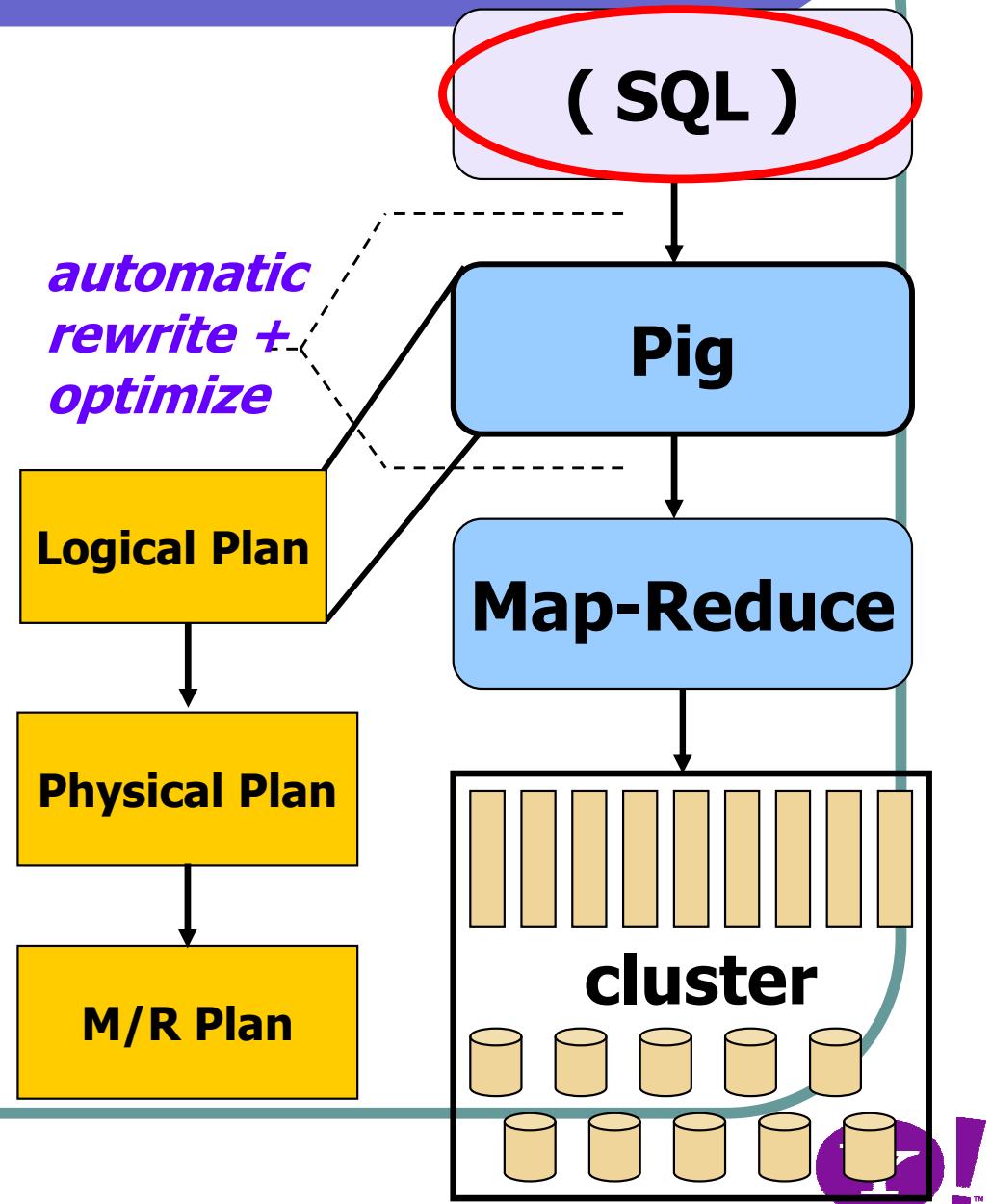
Srtd  = ORDER Smmd BY clicks;
Top100 = LIMIT Srtd 100;

STORE Top100 INTO 'top100sites.out';
```



# Pig script → Map/Reduce

- 不需懂底下 Map-Reduce 運作
- Pig 幫忙翻譯



# Why Pig?

- 容易學
- 開發快
- 一目瞭然

# Why Pig?

```

import java.io.IOException;
import java.util.ArrayList;
import java.util.Iterator;
import java.util.List;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.*;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobControl;
import org.apache.hadoop.mapred.KeyValueTextInputFormat;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.RecordReader;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.jobcontrol.Job;
import org.apache.hadoop.mapred.lib.IdentityMapper;
public class MRExample {
    public static class LoadPages extends MapReduceBase
        implements Mapper<LongWritable, Text, Text> {
        public void map(LongWritable k, Text val,
                       OutputCollector<Text, Text> oc,
                       Reporter reporter) throws IOException {
            // Pull the URL from the line
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String key = line.substring(0, firstComma);
            String value = line.substring(firstComma + 1);
            Text outkey = new Text(key);
            Text outVal = new Text(value);
            // Prepend an index to the value so we know which file
            // it came from.
            outVal.setIndex(1);
            oc.collect(outkey, outVal);
        }
    }
    public static class LoadAndFilterUsers extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text> {
        public void map(LongWritable k, Text val,
                       OutputCollector<Text, Text> oc,
                       Reporter reporter) throws IOException {
            // Pull the URL from the line
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String key = line.substring(0, firstComma);
            String value = line.substring(firstComma + 1);
            if (value.charAt(0) == '1')
                first.add(value.substring(1));
            else second.addValue.substring(1));
            reporter.setStatus("OK");
        }
        // Do the cross product and collect the values
        public void reduce(Text key, Iterable<Text> values,
                           OutputCollector<Text, Text> oc,
                           Reporter reporter) throws IOException {
            String outval = key + ',' + s1 + ',' + s2;
            oc.collect(null, new Text(outval));
            reporter.setStatus("OK");
        }
    }
}
public static void main(String[] args) {
    JobConf lp = new JobConf(MRExample.class);
    lp.setJobName("Load Pages");
    lp.setInputFormat(TextInputFormat.class);
    lp.setOutputKeyClass(Text.class);
    lp.setOutputValueClass(Text.class);
    lp.setMapperClass(LoadPages.class);
    lp.setInputFormat(FileInputFormat.class);
    Path("/user/gates/pages"));
    FileOutputFormat.setOutputPath(lp,
        new Path("user/gates/tmp/"));
    lp.setReducerClass(Reduce.class);
    Job loadPages = new Job(lp);

    JobConf lfu = new JobConf(MRExample.class);
    lfu.setJobName("Load and Filter");
    lfu.setInputFormat(TextInputFormat.class);
    lfu.setOutputKeyClass(Text.class);
    lfu.setOutputValueClass(Text.class);
    lfu.setMapperClass(LoadAndFilterUsers.class);
    lfu.setInputFormat(FileInputFormat.class);
    Path("/user/gates/users"));
    FileOutputFormat.setOutputPath(lfu,
        new Path("user/gates/tmp/"));
    lfu.setNumReduceTasks(0);
    Job loadUsers = new Job(lfu);

    JobConf join = new JobConf(MRExample.class);
    join.setJobName("Join Users and Pages");
    join.setInputFormat(KeyValueTextInputFormat.class);
    join.setOutputKeyClass(Text.class);
    join.setOutputValueClass(Text.class);
    join.setMapperClass(IdentityMapper.class);
    join.setReducerClass(Join.class);
    join.setInputFormat(FileInputFormat.class);
    Path("/user/gates/tmp/indexed_pages"));
    join.setOutputFormat(TextOutputFormat.class);
    join.setMapperClass(LoadUrls.class);
    join.setCombinerClass(SequenceFileCombiner.class);
    join.setReducerClass(ReduceUrls.class);
    join.setInputFormat(FileInputFormat.class);
    Path("/user/gates/tmp/joined"));
    join.setNumReduceTasks(50);
    Job joinJob = new Job(join);
    joinJob.addDependingJob(loadPages);
    joinJob.addDependingJob(loadUsers);

    JobConf group = new JobConf(MRExample.class);
    group.setJobName("Group URLs");
    group.setInputFormat(KeyValueTextInputFormat.class);
    group.setOutputKeyClass(Text.class);
    group.setOutputValueClass(LongWritable.class);
    group.setOutputFormat(TextOutputFormat.class);
    group.setMapperClass(LoadClicks.class);
    group.setCombinerClass(ReduceUrls.class);
    group.setReducerClass(ReduceUrls.class);
    group.setInputFormat(FileInputFormat.class);
    Path("/user/gates/tmp/joined"));
    FileOutputFormat.setOutputPath(group,
        Path("/user/gates/tmp/grouped"));
    group.setNumReduceTasks(50);
    Job groupJob = new Job(group);
    groupJob.addDependingJob(joinJob);

    JobConf top100 = new JobConf(MRExample.class);
    top100.setJobName("Top 100 sites");
    top100.setInputFormat(SequenceFileInputFormat.class);
    top100.setOutputKeyClass(LongWritable.class);
    top100.setOutputValueClass(Text.class);
    top100.setOutputFormat(SequenceFileOutputFormat.class);
    top100.setMapperClass(LoadClicks.class);
    top100.setCombinerClass(LimitClicks.class);
    top100.setReducerClass(LimitClicks.class);
    top100.setInputFormat(FileInputFormat.class);
    Path("/user/gates/tmp/grouped"));
    FileOutputFormat.setOutputPath(top100,
        Path("/user/gates/top100sitesforusers1"));
    top100.setNumReduceTasks(1);
    Job limit = new Job(top100);
    limit.addDependingJob(groupJob);

    JobControl jc = new JobControl();
    jc.addJob(loadPages);
    jc.addJob(loadUsers);
    jc.addJob(joinJob);
    jc.addJob(groupJob);
    jc.addJob(limit);
    jc.waitForCompletion();
}

```

**Users = LOAD 'users' AS (name, age);**

**Filtrd = FILTER Users by age >= 18 and age <= 25;**

**Pages = LOAD 'pages' AS (user, url);**

**Jnd = JOIN Filtrd BY name, Pages BY user;**

**Grpd = GROUP Jnd by url;**

**Smmtd = FOREACH Grpd GENERATE group, COUNT(Jnd) AS clicks;**

**Srted = ORDER Smmtd BY clicks;**

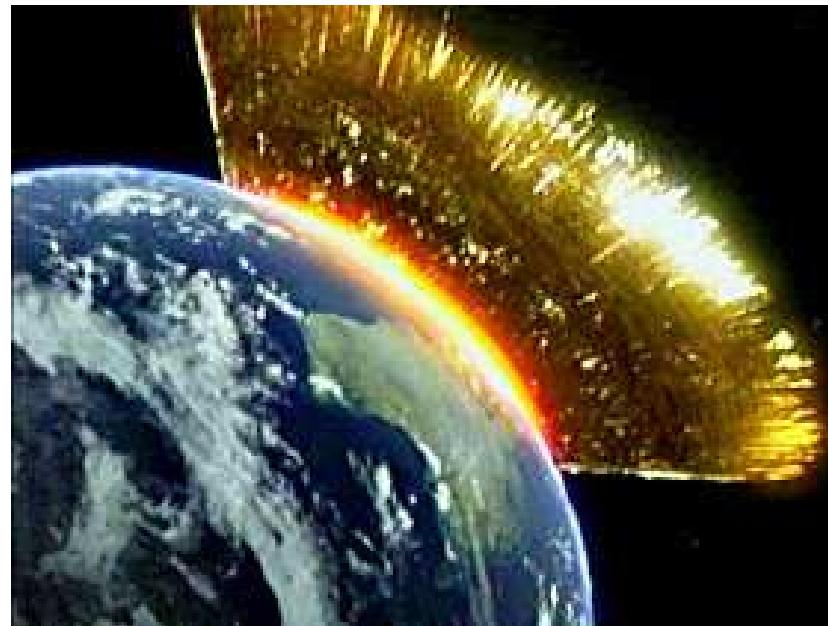
**Top100 = LIMIT Srted 100;**

**STORE Top100 INTO 'top100sites';**



# Why (NOT) Pig?

- 不是史上究極霹靂大無敵武器
  - Focus: aggregation, filter, join, ...
- 另一種做分散運算工作的方式

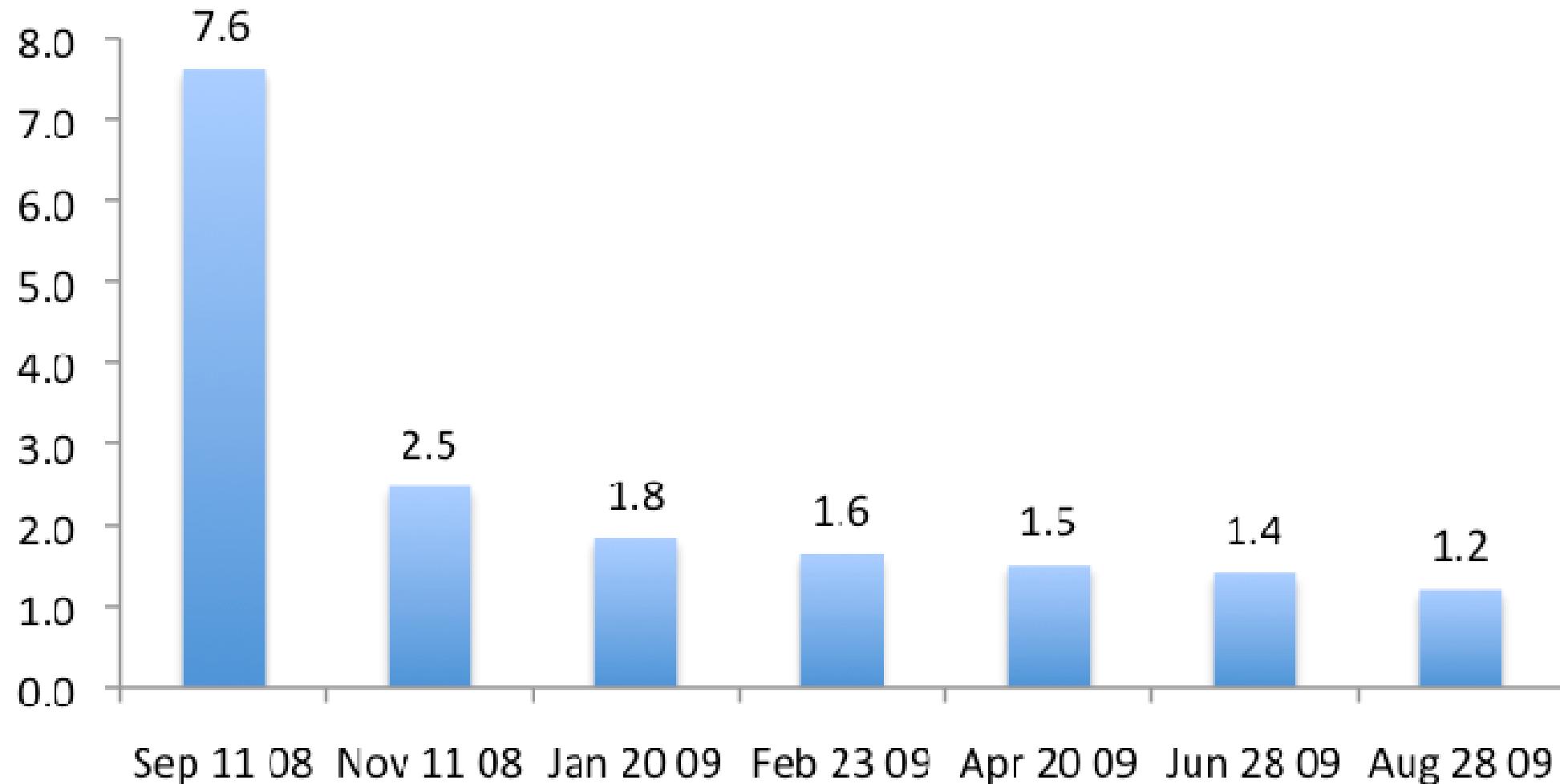


# Sweet spot between SQL – M/R

	SQL	Pig	Map-Reduce
<i>Programming style</i>	Large blocks of declarative constraints	→	“Plug together pipes”
<i>Built-in data manipulations</i>	Group-by, Sort, Join, Filter, Aggregate, Top-k, etc...	←	Group-by, Sort
<i>Execution model</i>	Fancy; trust the query optimizer	→	Simple, transparent
<i>Opportunities for automatic optimization</i>	Many	←	Few (logic buried in map() and reduce())
<i>Data Schema</i>	Must be known at table creation	→	Not required, may be defined at runtime



# Pig Performance vs Map-Reduce



# Execution and Syntax

# Pig Example

- Show users aged 18-25

```
Users = LOAD 'users.txt'  
          USING PigStorage(' , ') AS (name, age);  
Fltrd = FILTER Users  
          BY age >= 18 AND age <= 25;  
Names = FOREACH Fltrd GENERATE name;  
  
STORE Names INTO 'names.out';
```



# How to execute

- Local:

- `pig -x local foo.pig`

- Hadoop (HDFS):

- `pig foo.pig`
  - `pig -Dmapred.job.queue.name=xxx foo.pig`
    - `hadoop queue -showacls`



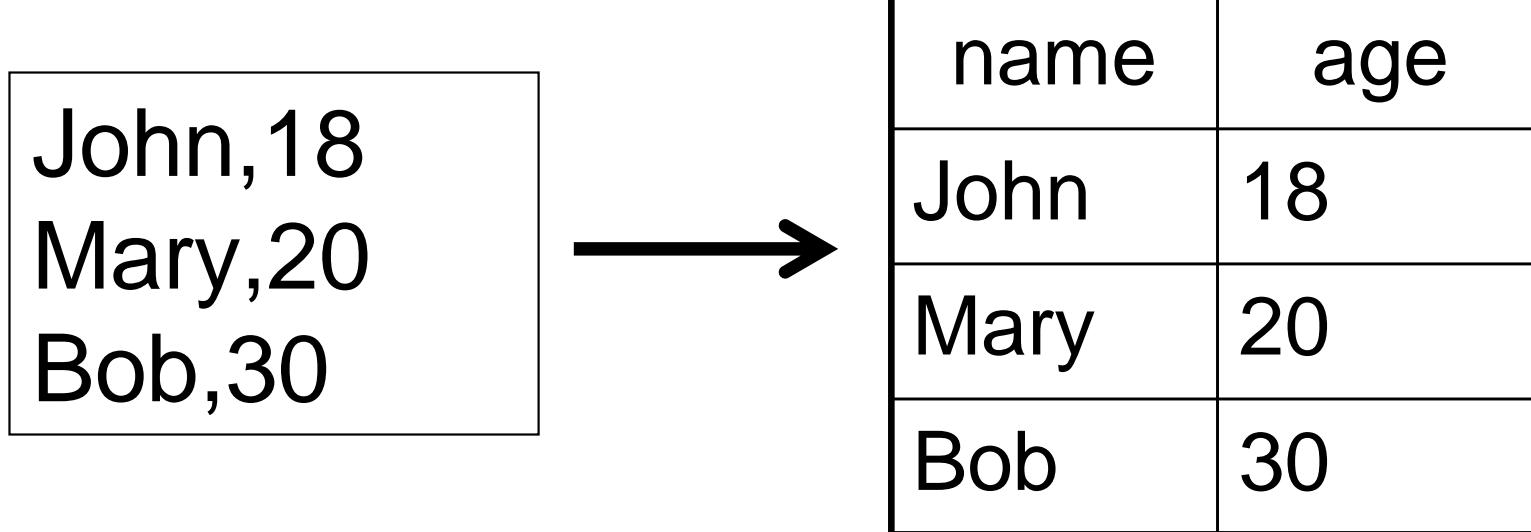
# How to execute

- Interactive pig shell
  - \$ pig
  - grunt> \_

# Load Data

```
Users = LOAD 'users.txt'  
          USING PigStorage(',') AS (name, age);
```

- LOAD ... AS ...
- PigStorage(',') to specify separator



# Filter

```
Fltrd = FILTER Users  
    BY age >= 18 AND age <= 25;
```

- FILTER ... BY ...
  - constraints can be composite

name	age
John	18
Mary	20
Bob	30



name	age
John	18
Mary	20

# Generate / Project

```
Names = FOREACH Fltrd GENERATE name;
```

- FOREACH ... GENERATE

name	age
John	18
Mary	20



name
John
Mary



# Store Data

```
STORE Names INTO 'names.out' ;
```

- STORE ... INTO ...
  - PigStorage(',') to specify separator if multiple fields



# Command - JOIN

```
Users = LOAD 'users' AS (name, age);  
Pages = LOAD 'pages' AS (user, url);  
Jnd   = JOIN Users BY name, Pages BY user;
```

name	age
John	18
Mary	20
Bob	30



user	url
John	yaho
Mary	goog
Bob	bing

name	age	user	url
John	18	John	yaho
Mary	20	Mary	goog
Bob	30	Bob	bing

# Command - GROUP

```
Grpd = GROUP Jnd by url;  
describe Grpd;
```

name	age	url
John	18	yhoo
Mary	20	goog
Dee	25	yhoo
Kim	40	bing
Bob	30	bing



yhoo	(John, 18, yhoo) (Dee, 25, yhoo)
goog	(Mary, 20, goog)
bing	(Kim, 40, bing) (Bob, 30, bing)

# Other Commands

- PARALLEL – controls #reducer
- ORDER – sort by a field
- COUNT – eval: count #elements
- COGROUUP – structured JOIN
- More at  
[http://hadoop.apache.org/pig/docs/r0.5.0/piglatin\\_reference.html](http://hadoop.apache.org/pig/docs/r0.5.0/piglatin_reference.html)



# Features

# Parameter Substitution

```
%default TYPE 'view'  
%declare ID '18987'  
A = load '/data/$DATE/$ID/$TYPE'
```

- \$ pig a.pig
- \$ pig -param DATE=20091009 a.pig
- \$ pig -param DATE=20091009 -param  
TYPE=click a.pig

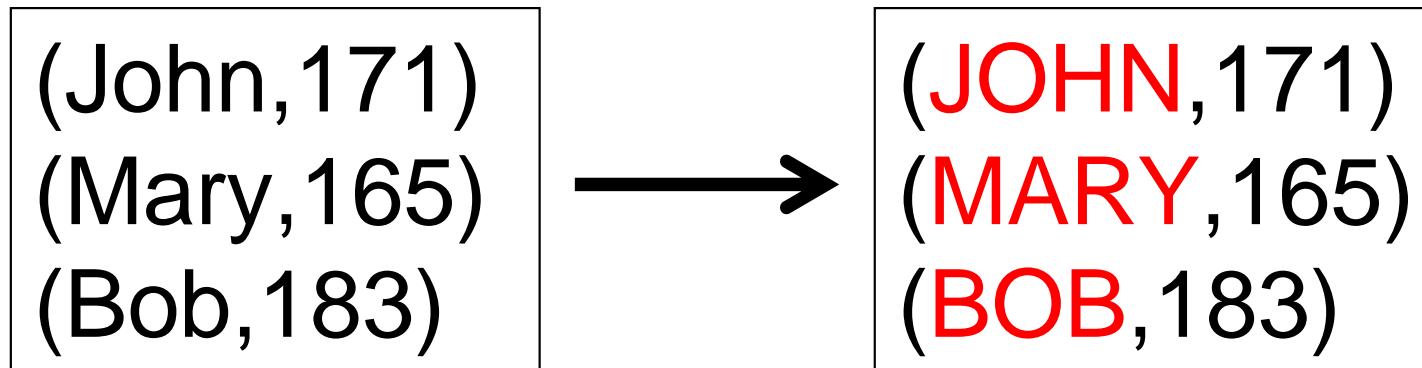


# RegEx Comparison

- itsyou = `FILTER urls by  
($0 MATCHES 'http://.*\\.yahoo\\.com.*')`
- MATCHES matches ‘whole’ string
  - ‘aaaa’ MATCHES ‘aaa.\*’ is true
  - ‘bbaaaa’ MATCHES ‘aaa.\*’ is false
- pattern syntax: `java.util.regex.Pattern`



# User-defined Function (UDF)



# UDF – user function part

```
package myudf;
import java.io.IOException;
import org.apache.pig.EvalFunc;
import org.apache.pig.data.Tuple;

public class UPPER extends EvalFunc<String>
{
    public String exec(Tuple in) throws IOException {
        if (in == null || in.size() == 0) return null;
        String str = (String)in.get(0);
        return str.toUpperCase();
    }
}
```



# UDF

- <http://hadoop.apache.org/pig/docs/r0.3.0/udf.html>
- <http://hadoop.apache.org/pig/javadoc/docs/api/>
- PiggyBank
  - Pig users UDF repo
  - <http://wiki.apache.org/pig/PiggyBank>



# Embedded in Java

```
/* create a pig server in the main class*/
{
    PigServer pigserver = new PigServer(args[0]);
    runMyQuery(pigServer, "/user/viraj/mydata.txt")
}

/* submit in function runMyQuery */

runMyQuery(PigServer pigServer, String inputFile) throws
IOException {
    pigServer.registerQuery("A = load '" + inputFile +
    "' as (f1,f2,f3);");
    pigServer.registerQuery("B = group A by f1;");
    pigServer.registerQuery("C = foreach B generate
flatten(group);");
    pigServer.store("C", "/user/viraj/myoutput");
}
```



# References

- FAQ
  - <http://wiki.apache.org/pig/FAQ>
- Documentation
  - <http://hadoop.apache.org/pig/docs/r0.5.0/>
- Talks & papers
  - <http://wiki.apache.org/pig/PigTalksPapers>
  - <http://www.cloudera.com/hadoop-training-pig-introduction>



# Questions?



# Backup slides

# Parameter Substitution

```
$ pig -param myparam=val foo.pig
```

```
B = filter A by ($0 eq '$myparam' )
```

- pig -dryrun produces processed script

```
B = filter A by ($0 eq 'val' )
```



# Parameter Substitution

- Params in file instead of command line
- \$ pig -param\_file myparams.txt a.pig

```
#myparams.txt  
DATE=20081009  
TYPE=clicks
```



# UDF – build user function

- javac
  - cp \$PIG\_HOME/lib/pig.jar**
  - sourcepath src
  - d classes
  - src/myudf/UPPER.java
- jar cf myudf.jar -C classes  
myudf/UPPER.class

# UDF – pig latin part

- **register** myudf.jar;
- B =  
foreach A generate  
**myudf.UPPER(name),height;**

# SQL vs. Pig Latin

<b><u>SQL</u></b>	<b><u>Pig</u></b>	<b><u>Example</u></b>
<b>From table</b>	<b>Load file(s)</b>	<b>SQL:</b> from X; <b>Pig:</b> A = load 'mydata' using PigStorage('\t') as (col1, col2, col3);
<b>Select</b>	<b>Foreach ... generate</b>	<b>SQL:</b> select col1 + col2, col3 ... <b>Pig:</b> B = foreach A generate col1 + col2, col3;
<b>Where</b>	<b>Filter</b>	<b>SQL:</b> select col1 + col2, col3 from X where col2>2; <b>Pig:</b> C = filter B by col2 > '2';

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<b><u>SQL</u></b>	<b><u>Pig</u></b>	<b><u>Example</u></b>
<b>Group by</b>	<b>Group + foreach ... generate</b>	<b>SQL:</b> select col1, col2, sum(col3) from X group by col1, col2; <b>Pig:</b> D = group A by (col1, col2); E = foreach D generate flatten(group), SUM(A.col3);
<b>Having</b>	<b>Filter</b>	<b>SQL:</b> select col1, sum(col2) from X group by col1 having sum(col2) > 5; <b>Pig:</b> F = filter E by \$1 > '5';
<b>Order By</b>	<b>Order ... By</b>	<b>SQL:</b> select col1, sum(col2) from X group by col1 order by col1; <b>Pig:</b> H = ORDER E by \$0;

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
<b>Distinct</b>	<b>Distinct</b>	<b>SQL:</b> select distinct col1 from X; <b>Pig:</b> I = foreach A generate col1; J = distinct I;
<b>Distinct Agg</b>	<b>Distinct in foreach</b>	<b>SQL:</b> select col1, count (distinct col2) from X group by col1; <b>Pig:</b> K = foreach D { L = distinct A.col2; generate flatten(group), SUM(L); }

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
Join	Cogroup + flatten  (also shortcut: JOIN)	<p><b>SQL:</b> select A.col1, B.col3 from A join B using (col1);</p> <p><b>Pig:</b></p> <pre>A = load 'data1' using PigStorage('t') as (col1, col2); B = load 'data2' using PigStorage('t') as (col1, col3); C = cogroup A by col1 <b>inner</b>, B by col1 <b>inner</b>; D = foreach C generate flatten(A), flatten(B); E = foreach D generate A.col1, B.col3;</pre>

(adapted from Viraj's slide)



# Debug Tips

- Use small data and pig -x local
- LIMIT
  - A = LOAD 'data' AS (a1,a2,a3)
  - B = LIMIT A 3;
- DUMP , DESCRIBE



# FAQ

- <http://wiki.apache.org/pig/FAQ>
  - can assign #reducer
  - support regex
  - can use allocated HOD cluster



# pig.vim

- [http://www.vim.org/scripts/script.php?script\\_id=2186](http://www.vim.org/scripts/script.php?script_id=2186)

```
A = load 'data.txt' as (f1,f2,f3);  
dump A;  
B = foreach A generate f1,f3;  
dump B;  
store B into 'output.txt' using PigStorage('-');
```





# Crawlzilla - A Toolkit for Deploying Cluster Search Engine Quickly and Easily

**Shun-Fa Yang 、 Wei-Yu Chen 、 Wen-Chieh Kuo**  
**Free Software Lab. @ NCHC**

INVENSIVE 2011 May 23, 2011

**TAIWAN**

[www.nchc.org.tw](http://www.nchc.org.tw)  
National Applied  
Research Laboratories



# Outline

Background and Motivation

- Nutch
- Hadoop
- Search Engine Library

## Crawlzilla

- Feature
- System Implement
- Demo

Introduction

Performance

Future Works





# Introduction

- The Information Explosion
- Increase Filter Efficiency by Search Engines
- Intranet also need Search Engines
- Build Search Engines isn't very Easy
- Crawlzilla can help You!



# Outline

Background and Motivation

- Nutch
- Hadoop
- Search Engine Library

## Crawlzilla

- Feature
- System Implement
- Demo

Introduction

Performance

Future Works



**CRAWLZILLA**

# Background and Motivation

Search Engine workflow

Related Open Source Projects

Compare with Other Projects

# Search Engine workflow – Phase 1

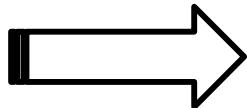
- Crawling the Web



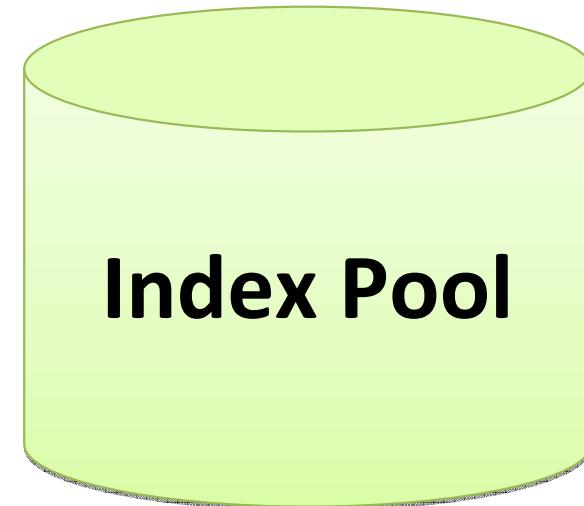
Crawler visits the web  
pages of the links

# Search Engine workflow – Phase 2

- Building the Index Pool

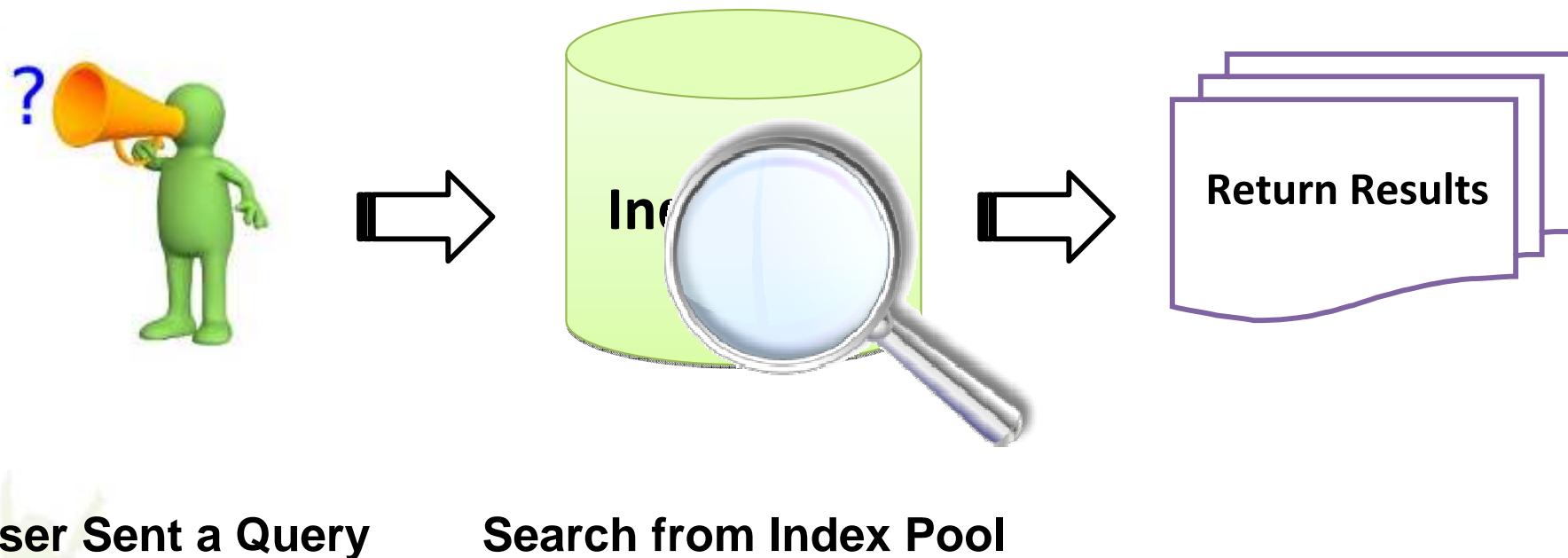


Parse Contents



# Search Engine workflow – Phase 3

- Serving Queries





# Background and Motivation

- **Related Open Source Projects**
  - Search Engine - Ntuch
  - Distributed Computing Platform – Hadoop
  - Search Engine Library – Lucene

# Background and Motivation

- If Build Search by Yourself ...
  - Setup Hadoop
  - Deploy System Configure Files
  - Debug Errors...
  - ...
  - ...
  - ...

# Compare with Other Projects

	Spidr	Larbin	Jcrawl	Nutch	Crawlzilla
Install	Rube Package Install	Gmake Compiler and Install	Java Compiler and Install	Deploy Configure Files	Provide Auto Installation
Crawl website pages	O	O	O	O	O
Parser Content	X	X	X	O	O
Cluster Computing	X	X	X	O	O
Interface	Command	Command	Command	Command	Web-UI
Support Chinese Segmentation	X	X	X	X	O

# Goal

- To Help Users to Build Search Engines Easily!
- To Help Users to Operate System Easily!
- Crawlzilla doesn't improve the algorithm of Nutch and Hadoop!
- Crawlzilla Provides Friendly Operating Interface and an Easy Way to Deploy Cluster Computing Environment!

# Outline

## Crawlzilla

- Feature
- System Implement
- Demo

## Background and Motivation

- Nutch
- Hadoop
- Search Engine Library

## Introduction

## Performance

## Future Works



# Crawlzilla Feature

- **Simply Install and Easy to Operate**
  - Customize user interface
- **More Powerful**
  - Support multiple search engines
- **More Search Engine Info.**
- **Developers to focus more**
  - Data mining tools



# Crawlzilla Architecture

**Web UI ( Crawlzilla Website + Search Engine)**

JSP + Servlet +  
JavaBean

Nutch

Lucene

**Crawlzilla System Management**

Tomcat

Hadoop

PC1

PC2

PC3

PC4

# System Implement

	JSP	Shell Script
Function	User UI	Admin and MIS UI
Security	Website Session	Crawler password with RSA Keys
Environment	Browser	Terminal with SSH –Client
Architecture	MVC	Module
Multi Language	i18n	Language parameters
	Default language is depend on O.S. Env.	

# Web Management

(Model 2)



## Setup and Drive Crawl

### Procedure

CrawlZilla Management Page

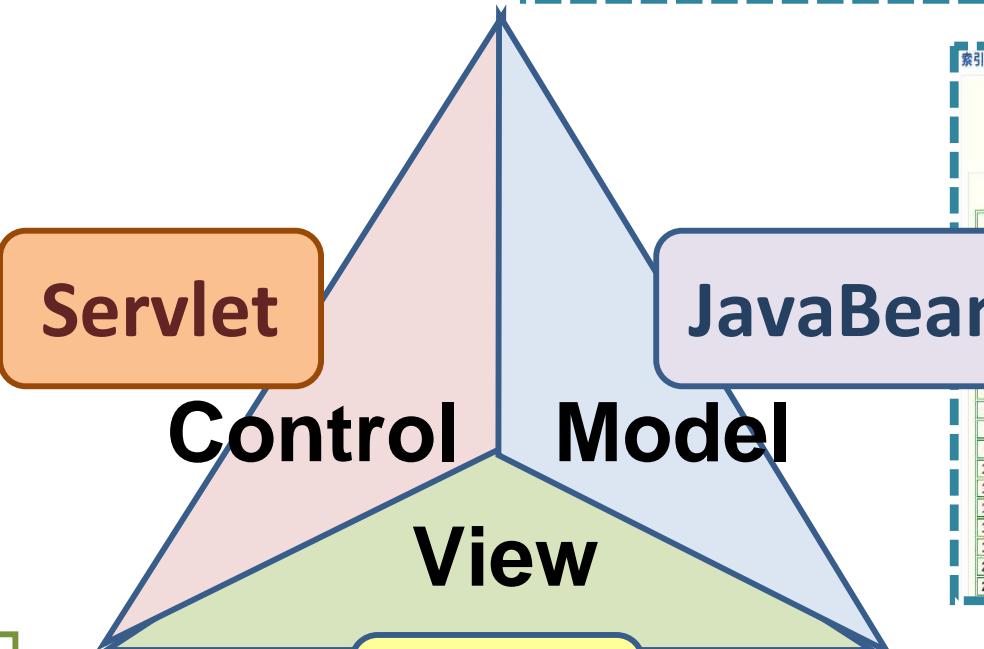
How To Use

1. Input Index Pool Name
2. Input URLs (see below example)  
http://www.nchc.org.tw  
http://www.google.com
3. Choose Depth, then Submit!

Index Pool Name: Index Pool Name

Input Crawl URLs: Input Crawl URLs

Crawl Depth Setup: Choose Crawl Depth



## System Status

CrawlZilla 網頁管理介面

系統狀態

索引庫狀態  
NCHC-2.3 crawling 0:0m51s [完成]

Jobtracker 工作排程器狀態 (New Window)

Running Jobs

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed
job_201009021521_0272	NORMAL	crawler	NCHC-2.3	50.00%	2	1	0.00%	1	0

Completed Jobs

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed
20100909134618									

## Session Certification

請輸入管理者密碼

\*\*\*\*\*

送出 重設

## Setup PW

這是你第一次登入 安全考量，預設的密碼不該被使用

原密碼為	*****
新設定的密碼	<input type="text"/>
確認新設定的密碼	<input type="text"/>

送出 重設

## Capture Lucene index pool

索引庫管理

索引庫名稱	建立時間	刪除索引庫	預覽	統計資料	嵌入後援引擎
nchc-en_3	2010-08-24 16:16:14	Delete	Preview	embed code	
nchc-tw_3	2010-08-24 15:22:48	Delete	Preview	embed code	

資料總覽

起始URL: http://www.nchc.org.tw/tw/

內容	引用次數	排序	內容	引用次數
www.nchc.org.tw	336	1	site:pccluster.nchc.org.tw	87
site:bioinfo.nchc.org.tw	66	3	site:nar.org.tw	57
site:edc.nchc.org.tw	53	5	site:service.nchc.org.tw	35
site:wanrc.nchc.org.tw	28	7	site:colife.nchc.org.tw	14
site:www.medicalgrid.org	13	9	site:volunteer.nchc.org.tw	9
site:www.spl.org.tw	7	13	site:nctwaren.net	7
site:ecgrid.nchc.org.tw	6	15	site:www.sipa.gov.tw	3
site:asp.104ehr.com.tw	3	17	site:viml.nchc.org.tw	3
site:www.ym.edu.tw	2	19	site:www.tnu.edu.tw	2
site:www.usc.edu.tw	2	21	site:www.ssvs.p.edu.tw	2
site:www.smlearning.org.tw	2	23	site:ecocam.nchc.org.tw	2

## i18N language setup

Home Crawl Manage System

Setup

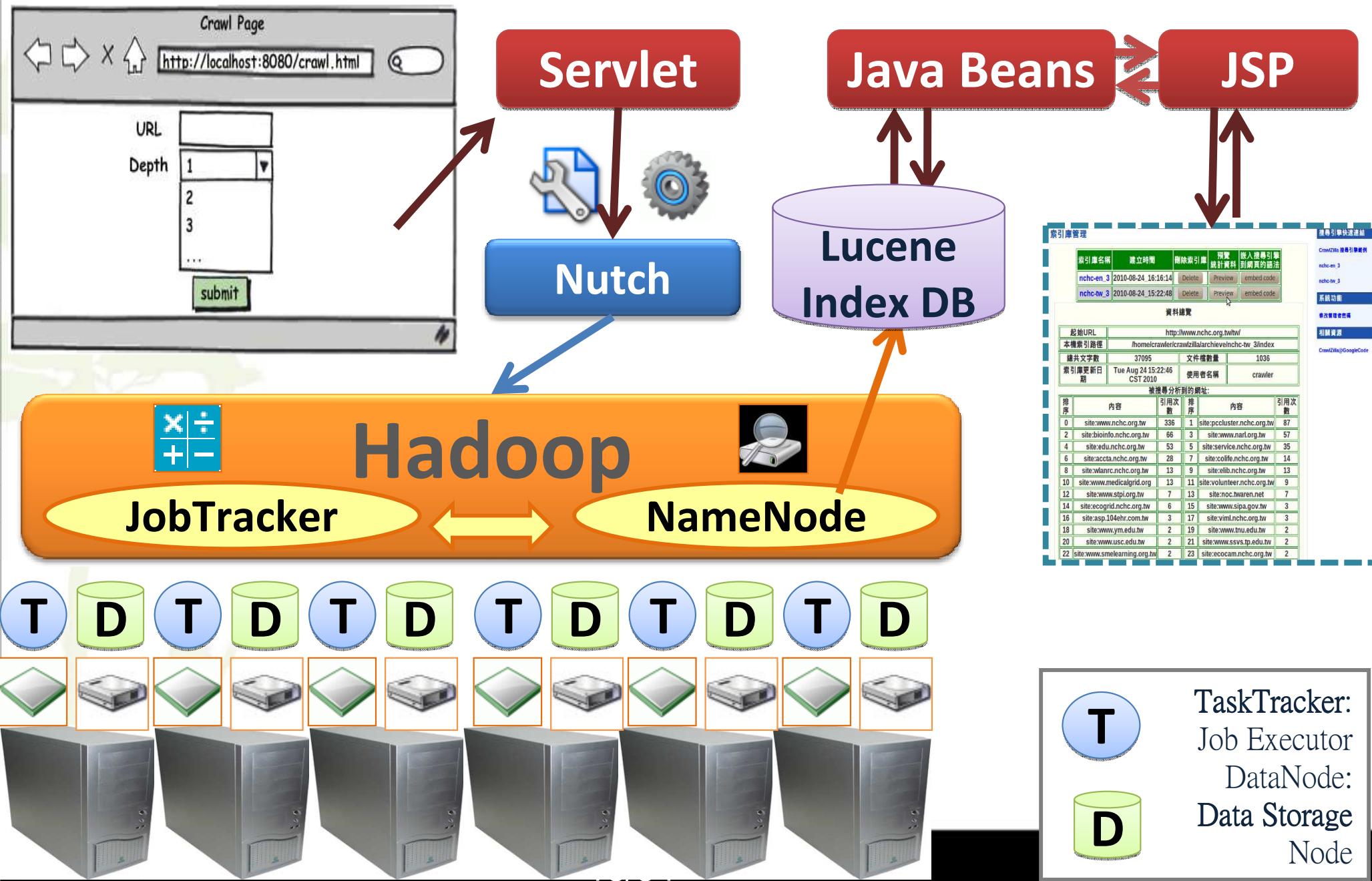
Engine Name: wuae

Admin Email: wuae@email

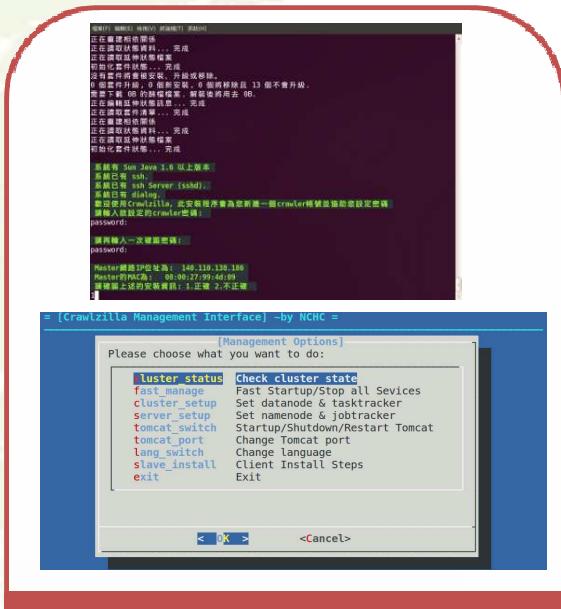
Choose Language: English

submit

# System Implement – Web Parser



# Friendly Interface!



**Admin**

CrawlZilla Management Page

Introduction of this Project

- Crawl Page : build up your search engine
- Index Pool Management : Setup and delete the result Index Pool
- System Status : Inspect your result Index Pool
- Admin Setup : Setup personality and multi-language
- Setup the password at this web page

Index Pool Name	Created Time	Crawling Depth	Crawling Time	Delete Index Pool	Preview Statistics Data	Re Crawl	code of ended search bar to web page
udn-3	2011-01-24 14:36:54	3	0m53m58s	<a href="#">Delete</a>	<a href="#">Preview</a>	<a href="#">ReCrawl</a>	<a href="#">ended code</a>

Data Overview

Initial Urls	http://udn.com/NEWS/mainpage.shtml
Local Index Path	/home/crawler/crawlzillarchive/udn-3/index
Total Words	89168
Total Files	4642
Index Pool Updated Time	Mon Jan 24 14:36:54 CST 2011
User Name	crawler

parsed URLs:

順序	內容	引用次數	順序	內容	引用次數
0	site:mag.udn.com	1159	1	site:udn.com	537
1	site:money.udn.com	401	3	site:travel.udn.com	316
4	site:stars.udn.com	313	5	site:video.udn.com	309
6	site:udn.gohappy.com.tw	244	7	site:blog.udn.com	180
8	site:digiveup.udn.com	158	9	site:pro.udnjob.com	129
10	site:learning.udn.com	123	11	site:bookmark.udn.com	120
12	site:udn.com	111	13	site:viva.udn.com	93
14	site:health.udnjobs.com	74	15	site:stock.udn.com	56
16	site:album.udn.com	49	17	site:www.udngroup.com	46
18	site:udn.magazine.com.tw	43	19	site:reporter.udn.com	40
20	site:co.udn.com	25	21	site:www.gohappy.com.tw	12

**MIS**

CrawlZilla Management Page

Cloud

ca | de | en | es | fr | hu | it | ja | cs | nl | pt | pl | st | sr | sv | th | zh | cn |

NCHC 國立應用科學研究中心  
National Center for High-Performance Computing  
Better HPC Better Living

Powered by 

**USER**

# Live Demo I

## Crawlzilla Install

- (1) Master Install
- (2) Cluster Slave Install



# Live Video Demo

- Master Install ([Demo Video also @ YouTube](#))





# Live Video Demo

- Slave Install ([Demo Video also @ YouTube](#))



# Live Demo II

## Dialog Management

# Live Demo III

## Web Management

- (1) Crawl Setup
- (2) Search Engine Index Pool
- (3) Search it!

# Outline

Background and Motivation

- Nutch
- Hadoop
- Search Engine Library

## Crawlzilla

- Feature
- System Implement
- Demo

Introduction

Performance

Future Works

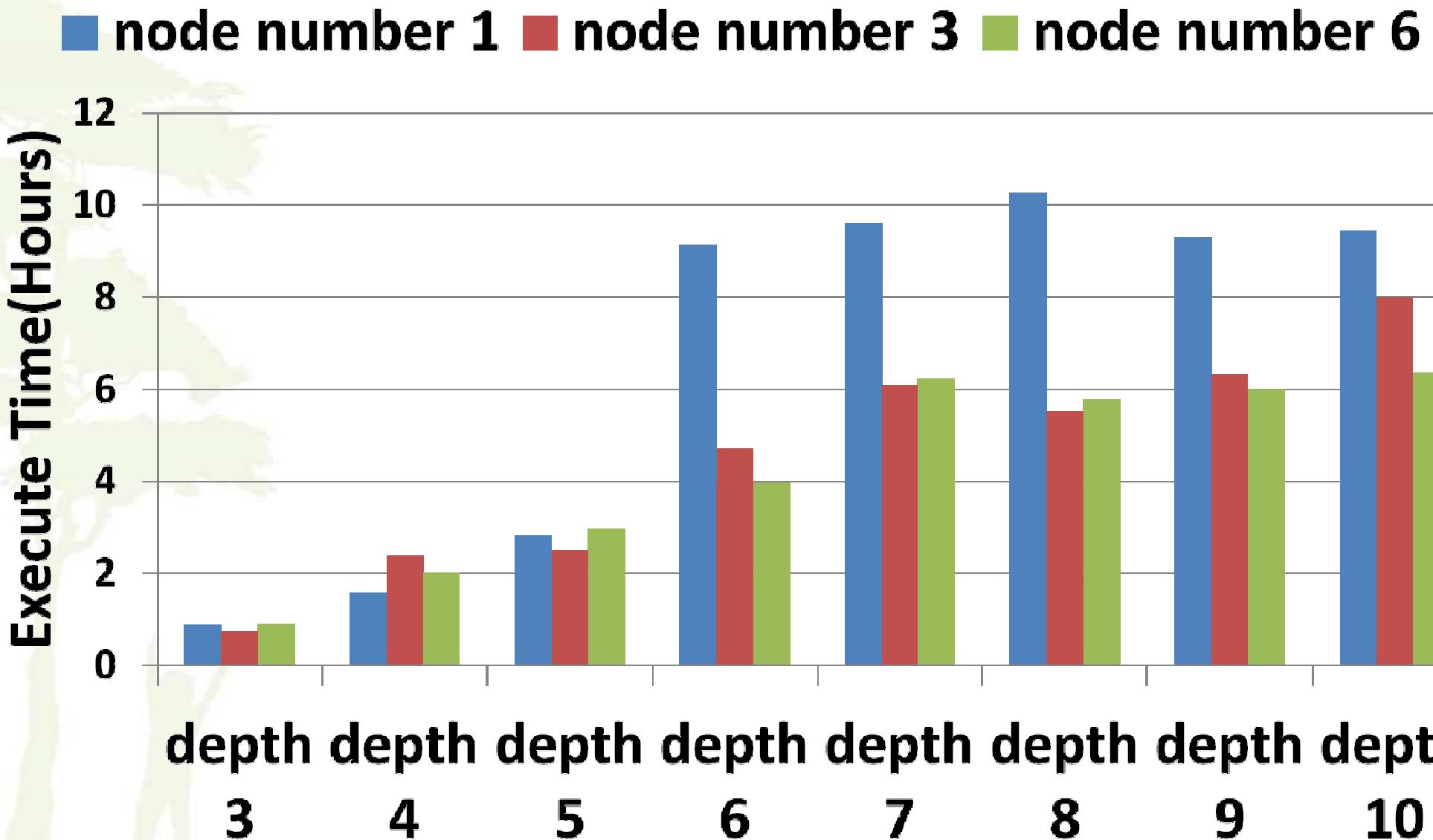


# Performance

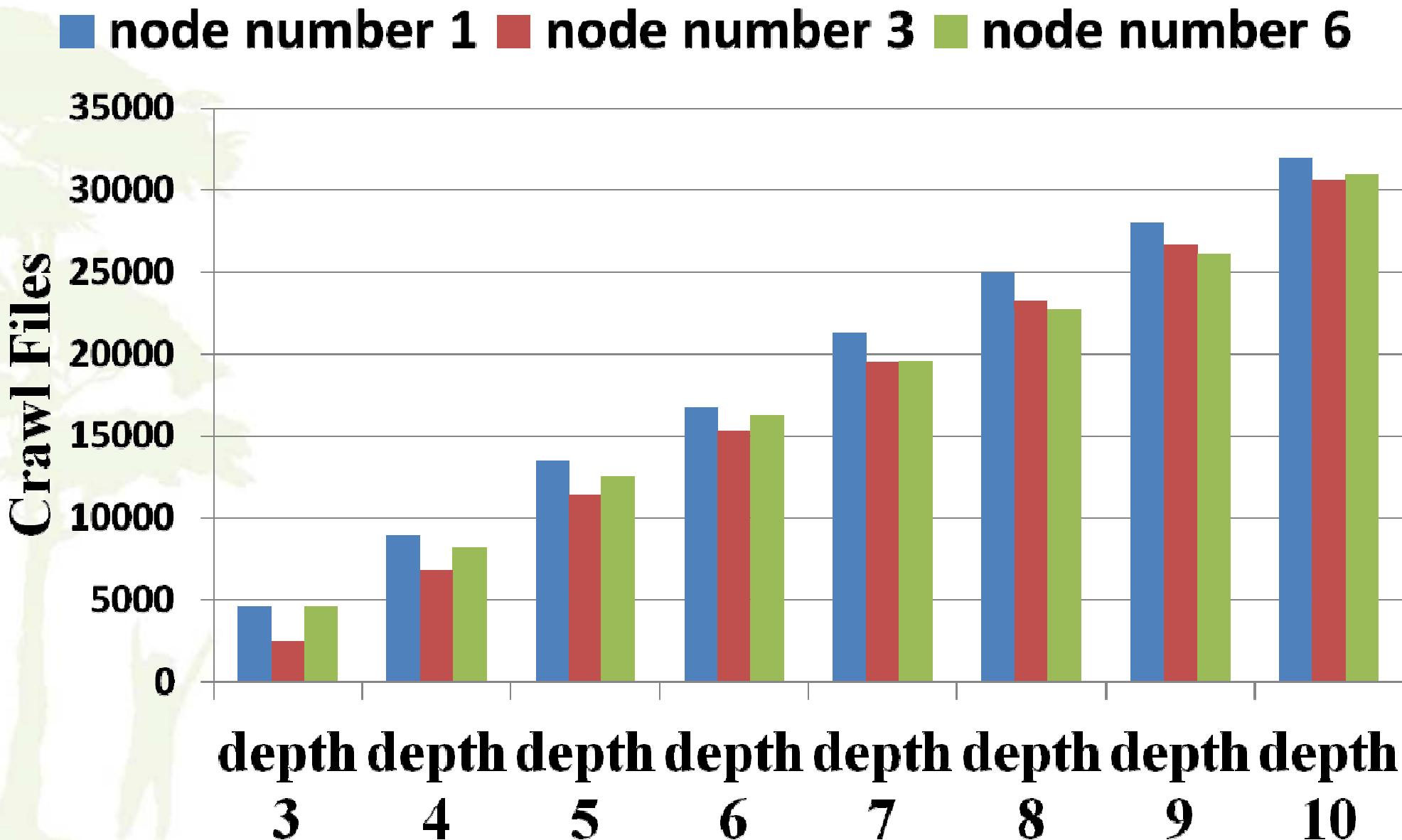
## Experiment Environment

- CPU
  - Intel(R) Core(TM)2 Quad CPU Q9550 2.83GHz
- Memroy
  - 8 GigaBytes
- Operation System
  - Ubuntu 10.04 Lucid(x86)
- Crawlzilla Version
  - 0.3.0-101116

# Execute Time

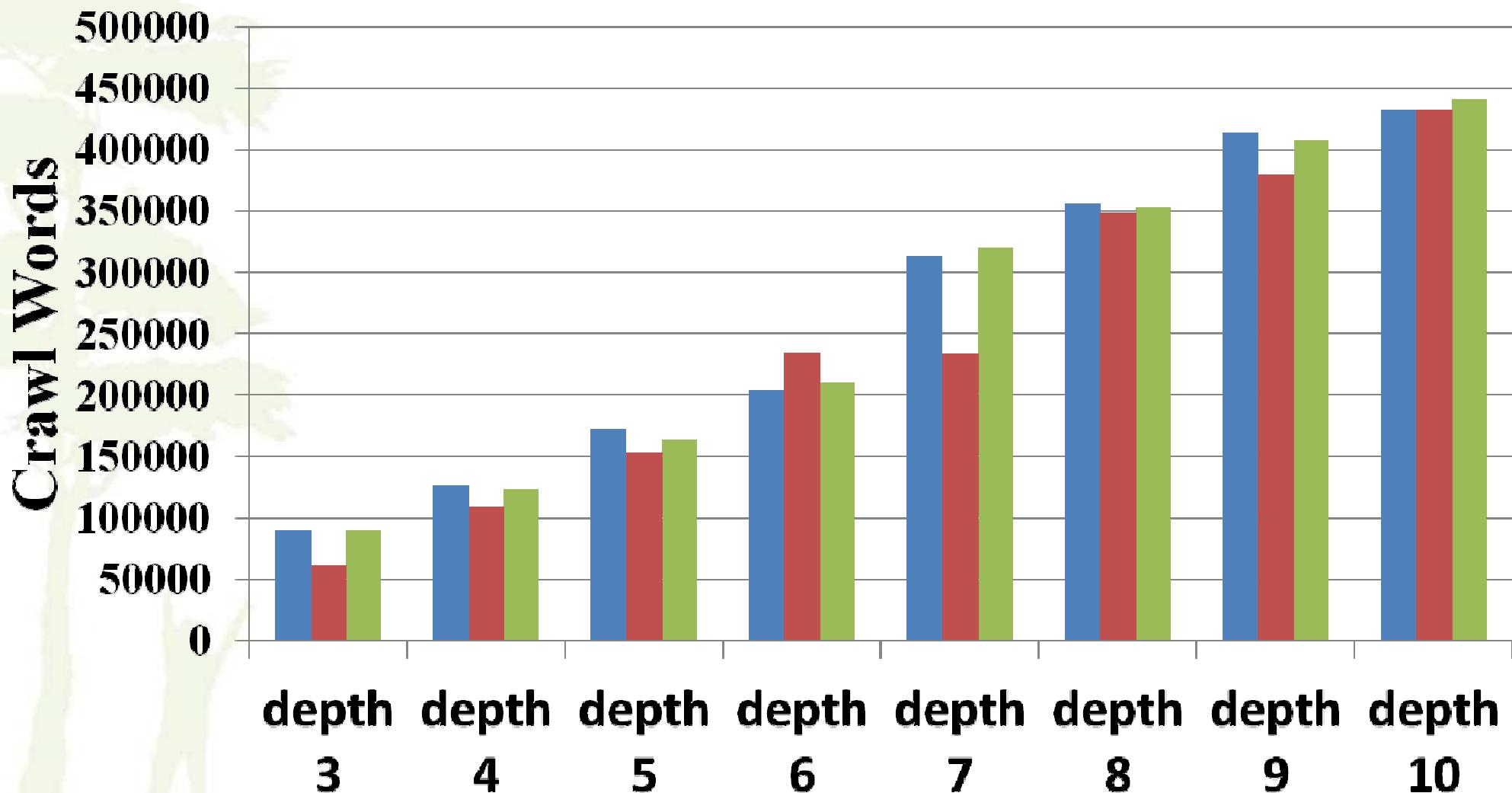


# Crawl Files



# Crawl Words

■ node number 1 ■ node number 3 ■ node number 6



# Outline

Background and Motivation

- Nutch
- Hadoop
- Search Engine Library

## Crawlzilla

- Feature
- System Implement
- Demo

Introduction

Performance

Future Works





# Future Works

- **New Version**

- Support Multi User
- Support Schedule
- Update the Kernel
- More Easily to deploy Slave Computing Nodes
- Now is testing!
- Release Day See <http://crawlzilla.info>

# Reference

- **J. Dean and S. Ghemawat, MapReduce: Simplified Data Processing on Large Clusters, In Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation - Volume 6, San Francisco, CA, December 06 - 08, 2004.**
- **S. Ghemawat, H. Gobioff and S. T. Leung, The Google File System, 19th ACM Symposium on Operating Systems Principles, Lake George, NY, October, 2003.**
- **The Apache Software Foundation, Nutch, available at:  
<http://nutch.apache.org/>, accessed 5 June 2010.**
- **The Apache Software Foundation, Hadoop, available at:  
<http://hadoop.apache.org/>, accessed 5 June 2010.**
- **The Apache Software Foundation, Lucene, available at:  
<http://lucene.apache.org/>, accessed 5 June 2010.**
- **Crawlzilla @ Google Code Project Hosting, available at:  
<http://code.google.com/p/crawlzilla/>, accessed 15 Sep 2010.**



# Enjoy your search engines!!!

## Start from Here!

- Crawlzilla @ Google Code Project Hosting (Tutorials in Chinese)
  - <http://code.google.com/p/crawlzilla/>
- Crawlzilla @ Source Forge (Tutorials in English)
  - <http://sourceforge.net/p/crawlzilla/home/>
- Crawlzilla User Group @ Google
  - <http://groups.google.com/group/crawlzilla-user>
- NCHC Cloud Computing Research Group
  - <http://trac.nchc.org.tw/cloud>

# Thank You!

## Q & A





# 運用自由軟體打造資安雲端分析平台

Building Network Security Cloud Analysis Platfrom using Open Source

**Yao-Tsung Wang**

[jazz@nchc.org.tw](mailto:jazz@nchc.org.tw)

**Wei-Yu Chen**

[wuae@nchc.org.tw](mailto:wuae@nchc.org.tw)



# 專家說：雲端每個環節都有安全問題



ZDNet Taiwan - 專家談雲端：每個環節都有安全問題 - 新聞



2010/08/10 19:50:02



專家談雲端：每個環節都有安全問題

ZDNet記者曠文濬／台北報導 雲端的安全問題不是無解，只是不管是雲端服務供應商或者想要建立私有雲的企業用戶，都必須考量到每個環節。

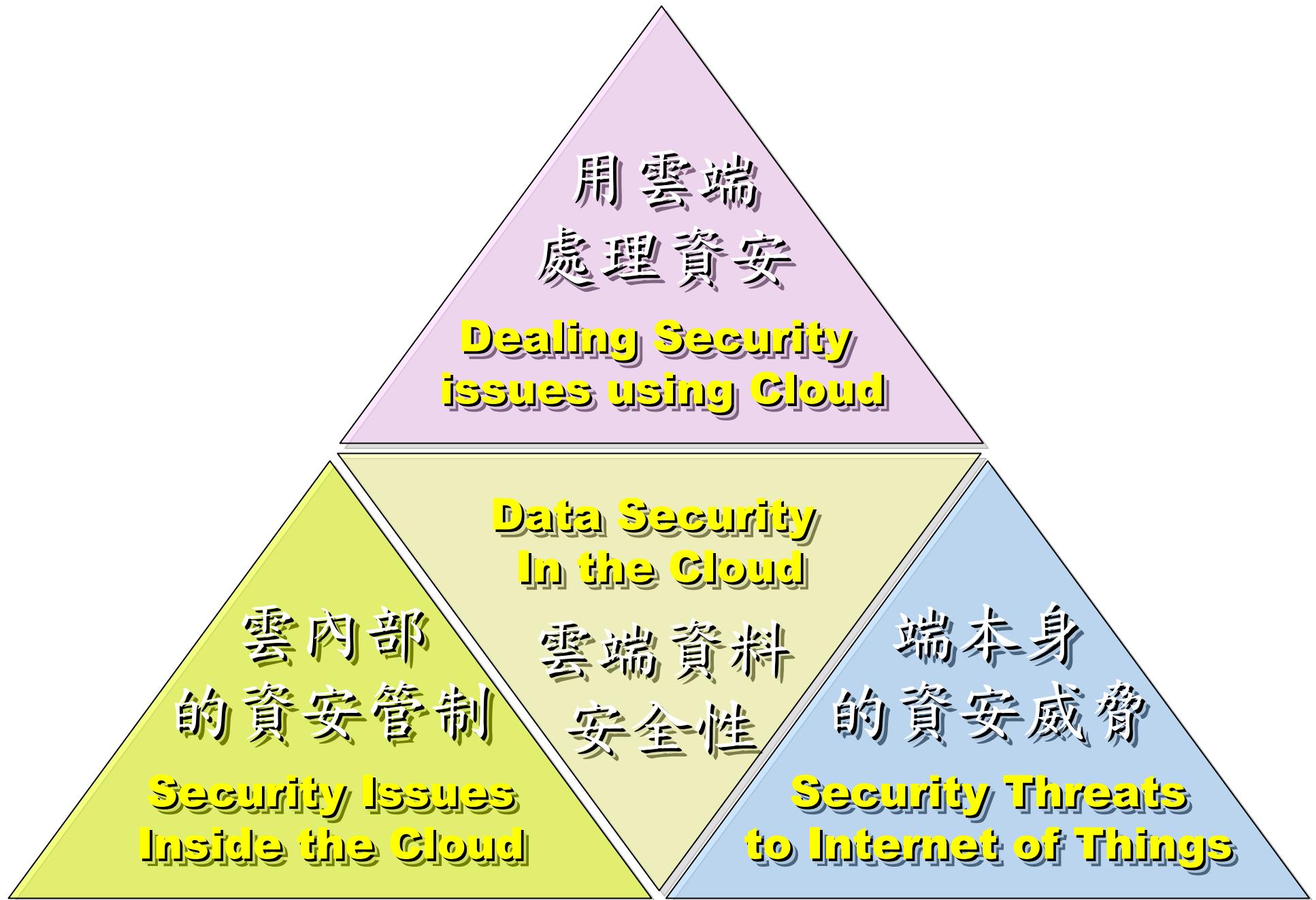
微軟亞太區全球技術支援中心專案經理、同時也是ZDNet專欄作家林宏嘉今（10）日在ZDNet舉行的IT Priorities圓桌論壇中表示，**雲端的安全議題涉及了IaaS、PaaS乃至於SaaS的每個層面**。當然有些問題是原本就存在：例如在討論到IaaS時，就涉及到了**機房的管理和硬體設備的可用性**等；但是講到PaaS時，企業用戶倘若要選擇開原碼的作業系統，必須考量到後續的**安全維護**；在SaaS的層次，企業用戶必須確保每一個分區（partition）的安全更新和**資料安全**。

目前正如火如荼建立台灣第一個校園私有雲的台大計算機及資訊網路中心主任孫雅麗則呼應道，Amazon的雲端服務證實了在Hypervisor層有駭客入侵，也就是意味著過去大家在討論如何防範**虛擬機器的資料安全**，但是威脅已經深化到了更下一層。這些問題都有待解決。

「有些問題甚至是來自於內部，舉例而言，MIS可能會把存在記憶體裡的資料倒出來，或者在Hypervisor層就植入了可以蒐集資料的程式，」孫雅麗說。

安全議題是目前台灣企業對雲端持保留態度的最大主因，這也是何以台灣的大型企業對於雲端的想法，還是抱持懷疑和反彈，畢竟中國公司，連同大陸的政府，都是云頂娛樂城「甘肅一派」。

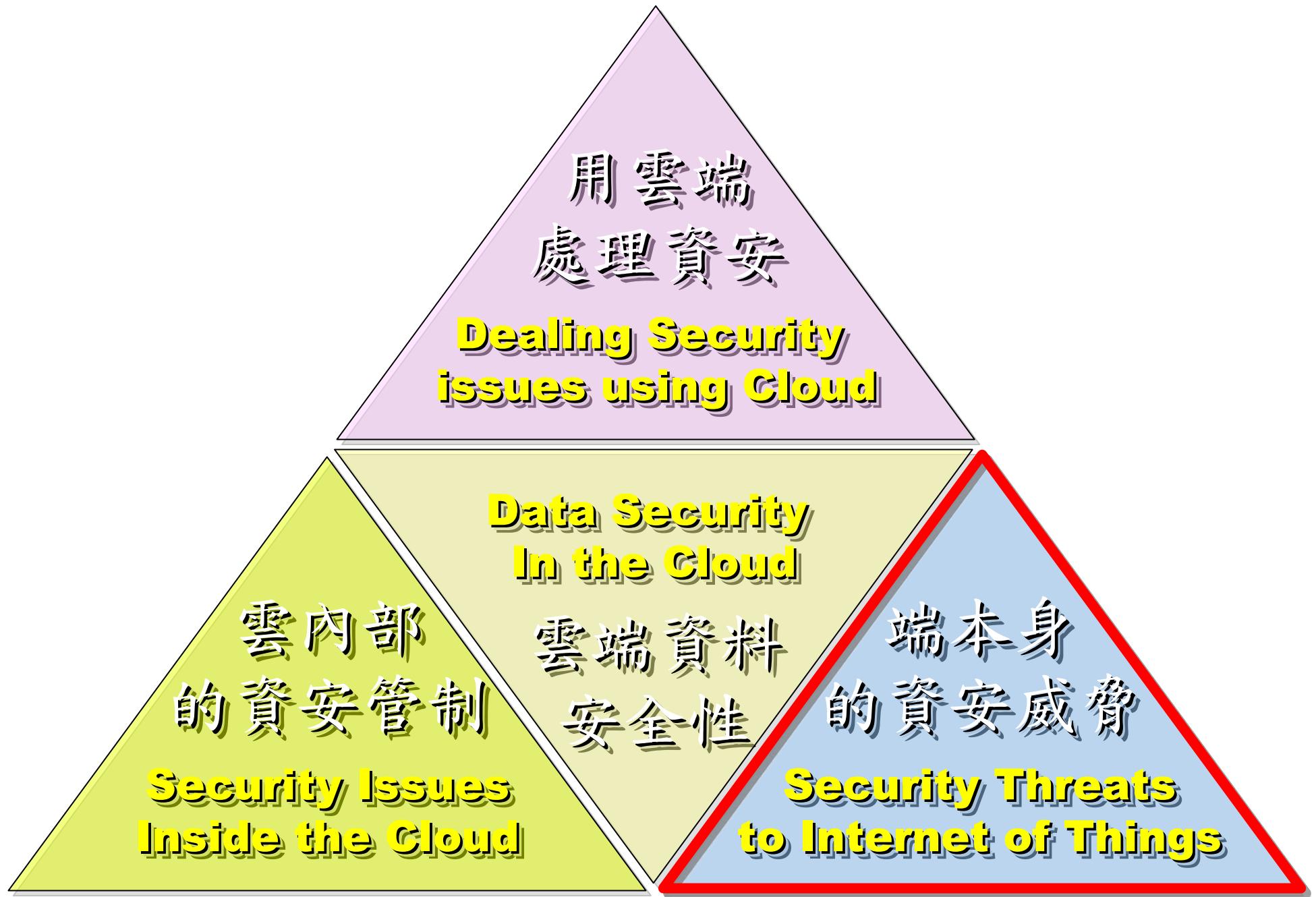
# 雲端資安的範疇



# 兩大研究方向：你該選「雲」還是「端」？



# 先來談談「端的安全」



以前你只有電腦需要防毒，現在 .....



端

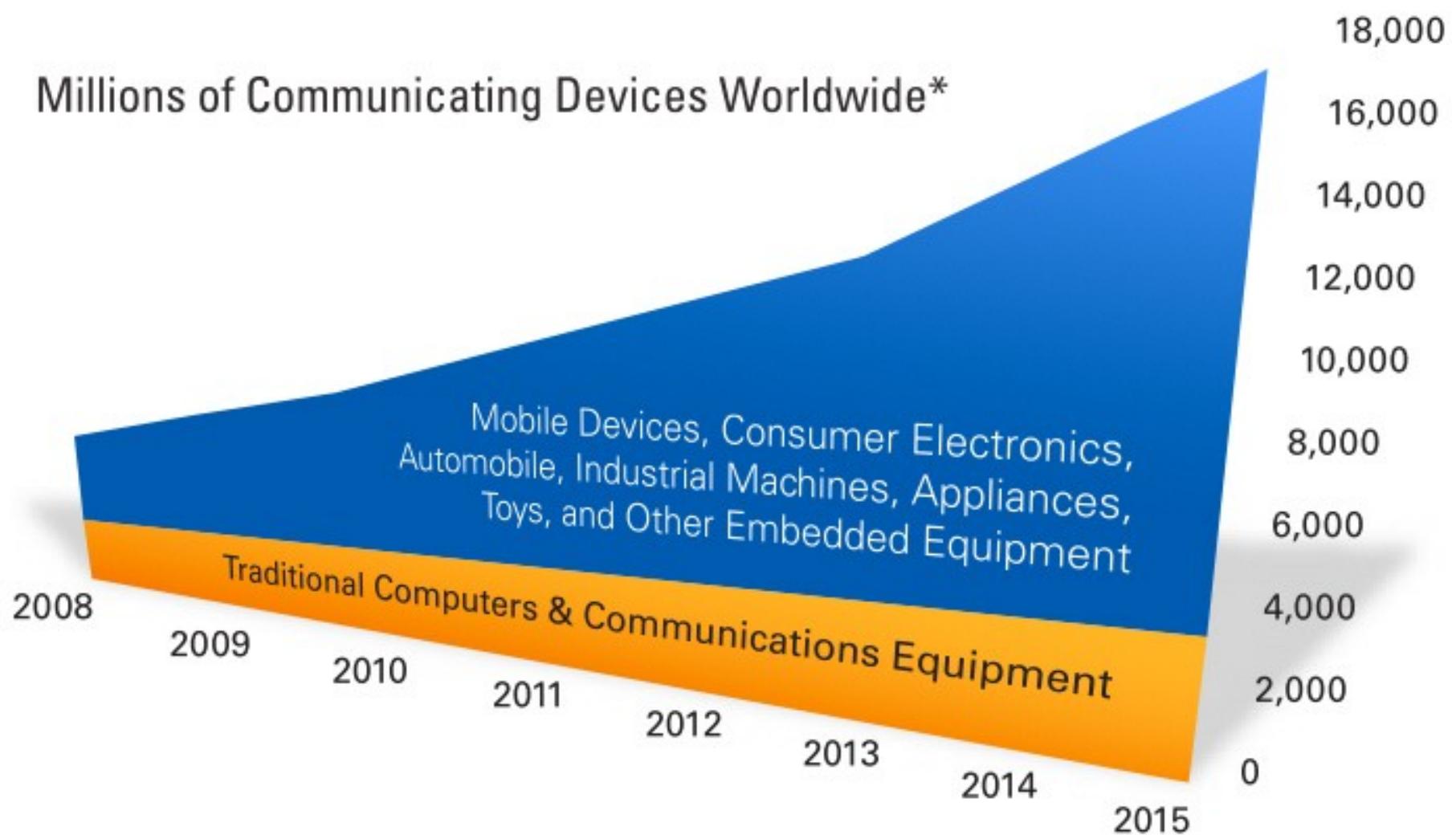


symbian  
OS



多元，中小廠  
Diversify ,  
SMB

# 全球連網裝置急速成長中



Source: IDC Device Base Model, 2009

\*Excludes voice- and SMS-only phones

# 物聯網的時代來臨

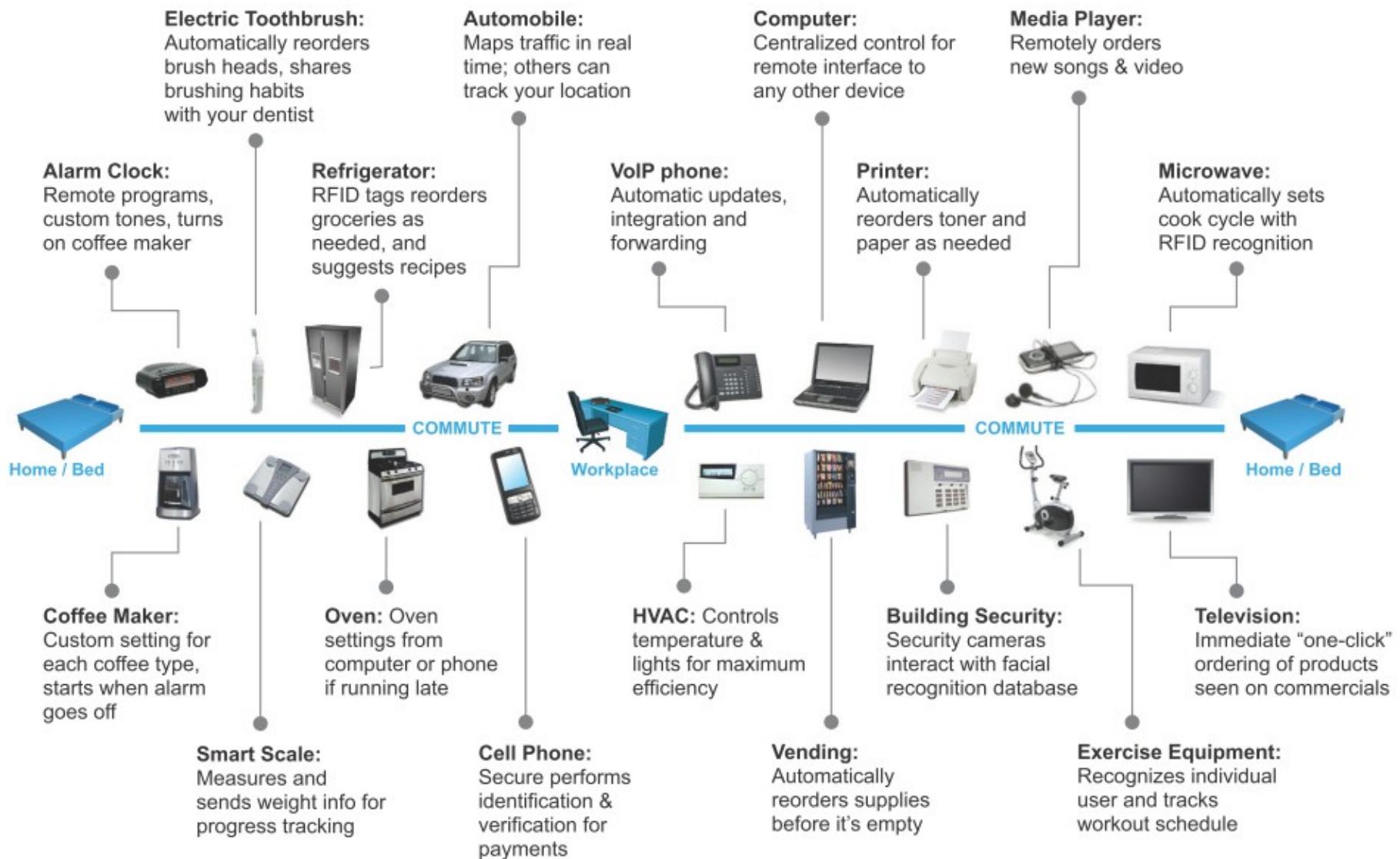
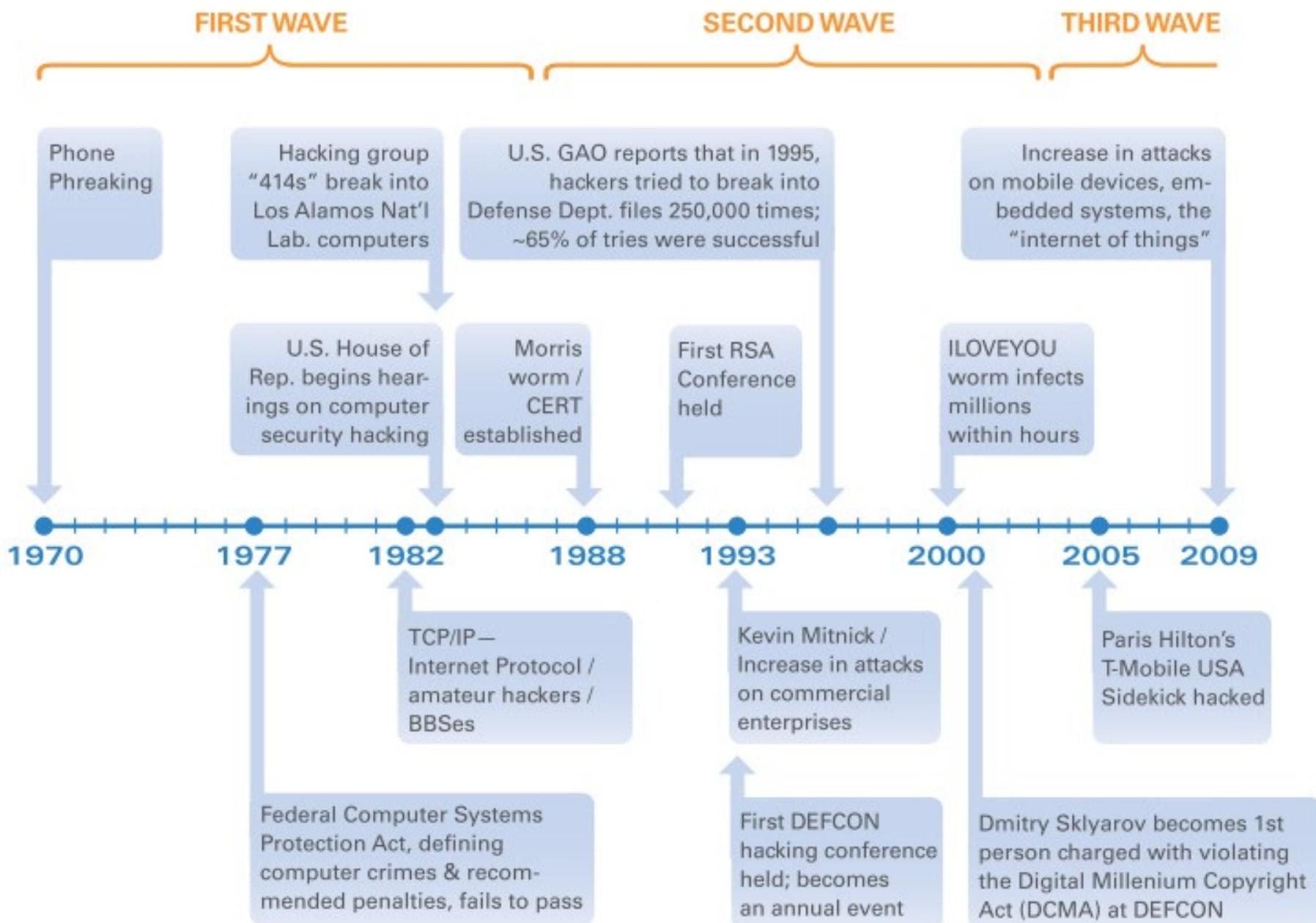


Figure 3. The Internet of Things

圖片來源：[Attacks on Mobile and Embedded Systems: Current Trends by Mocana](#)

# 第三波網路入侵對象將鎖定在『物聯網』



圖片來源：[Attacks on Mobile and Embedded Systems: Current Trends by Mocana](#)

# 針對行動裝置的各種資安問題與經驗

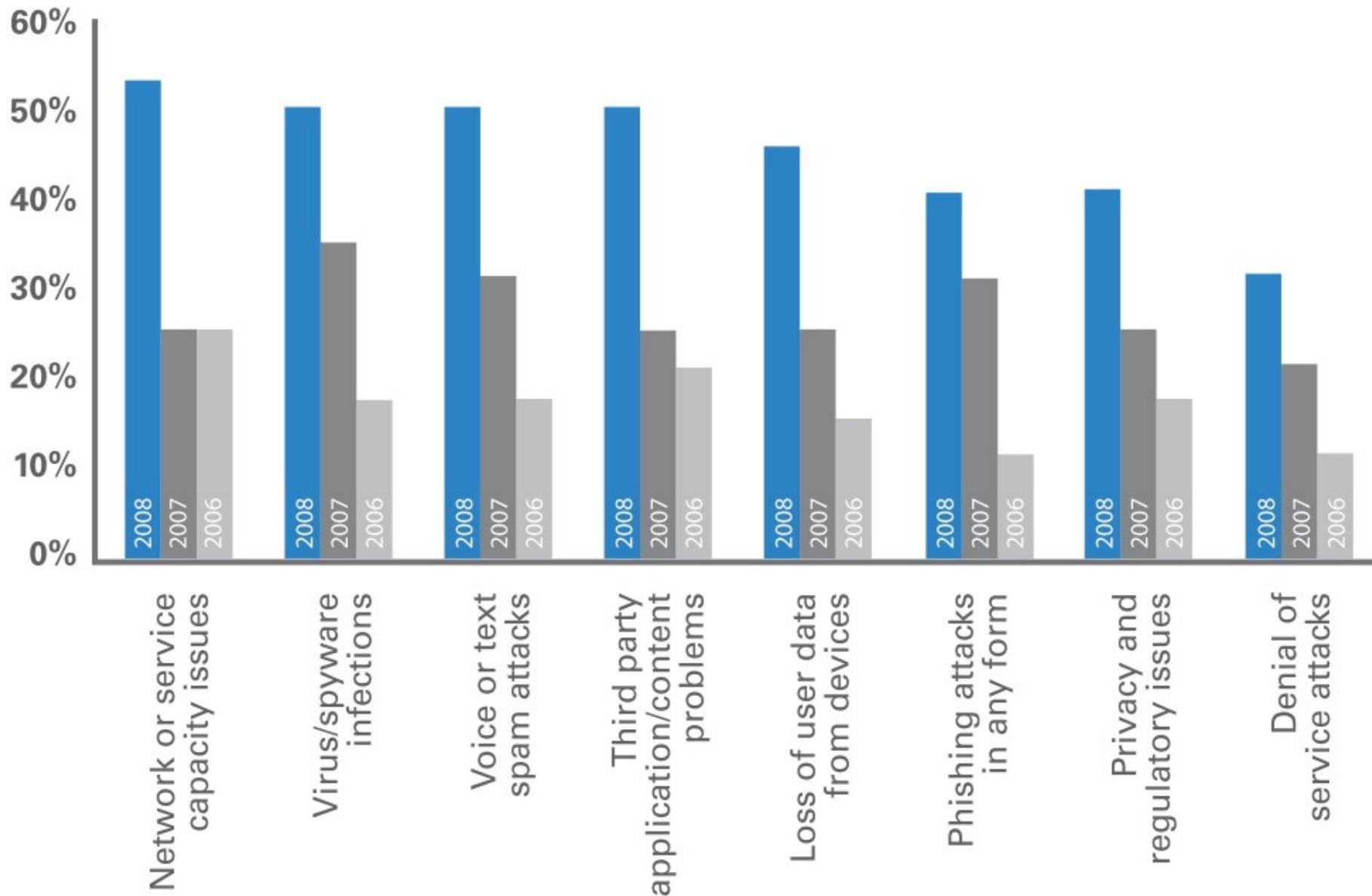
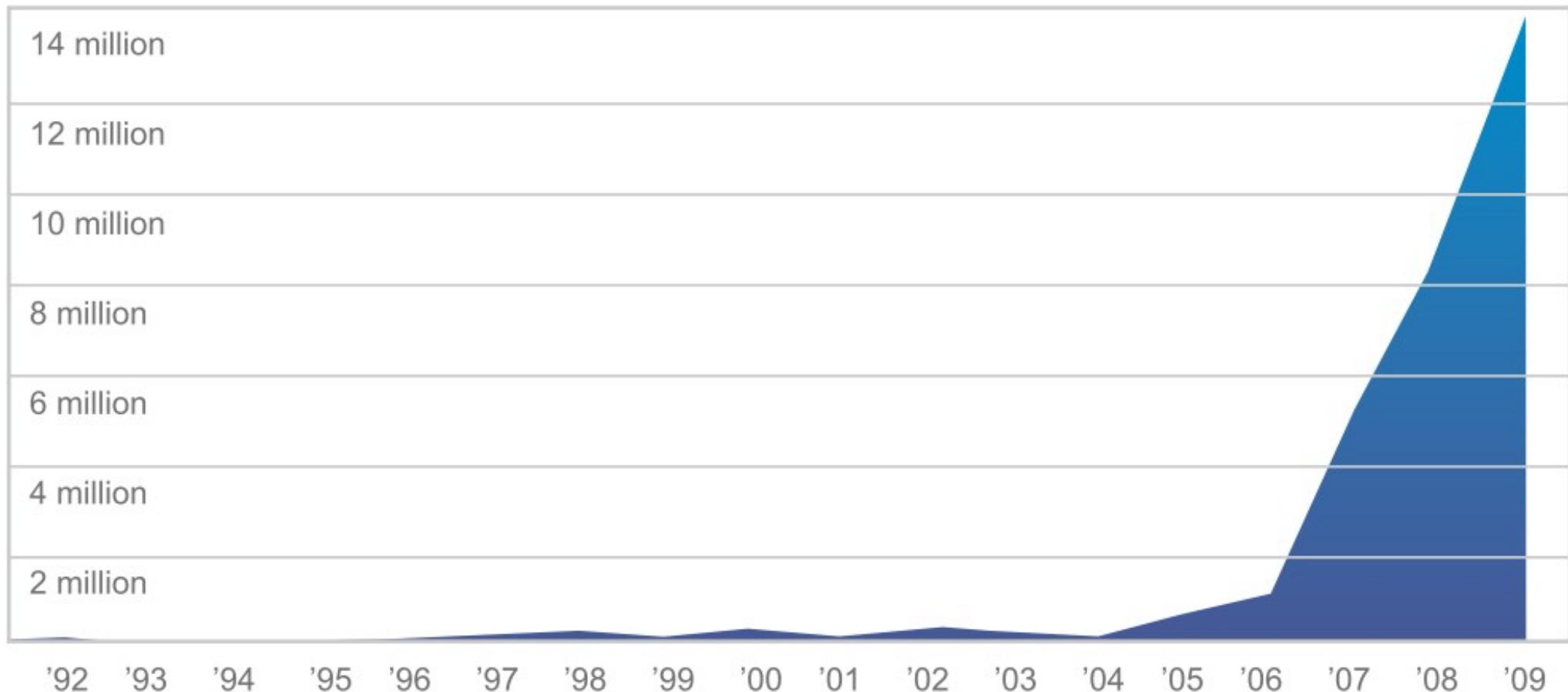


Figure 6. The increase in security issues experienced by mobile device users from 2006 to 2008; % of respondents. McAfee Mobile Security Report 2009

圖片來源：[Attacks on Mobile and Embedded Systems: Current Trends by Mocana](#)

# 網路惡意程式 (Malware) 逐年激增

Malware detected by year



Over 3,000 new “species” of PC malware are released onto the Internet every hour. Now that malware is setting its sights on Device platforms.

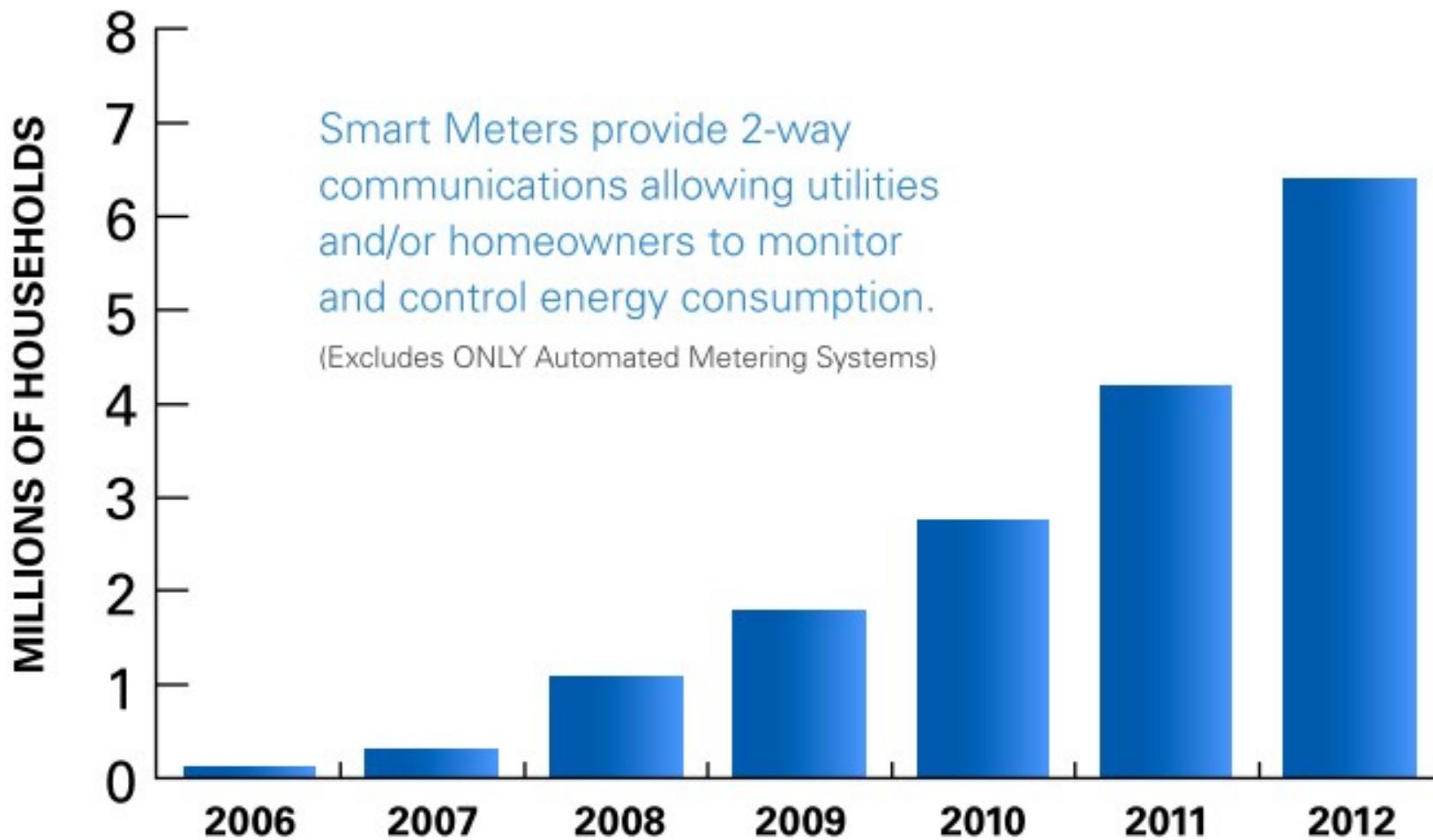
Source: AV LABS

圖片來源：

[U.S. Unprepared for Internet Device Flood: Unaddressed Security Problems & Talent Drought Threaten Long-Term Commercial, Government Interests](#) 11  
By: Kurt Stammberger, CISSP, Adrian Turner and Mat Small, Mocana With: Rich Nass, Sarah Friar, Goldman Sachs

# 如果你家的智慧電錶被入侵會怎樣？

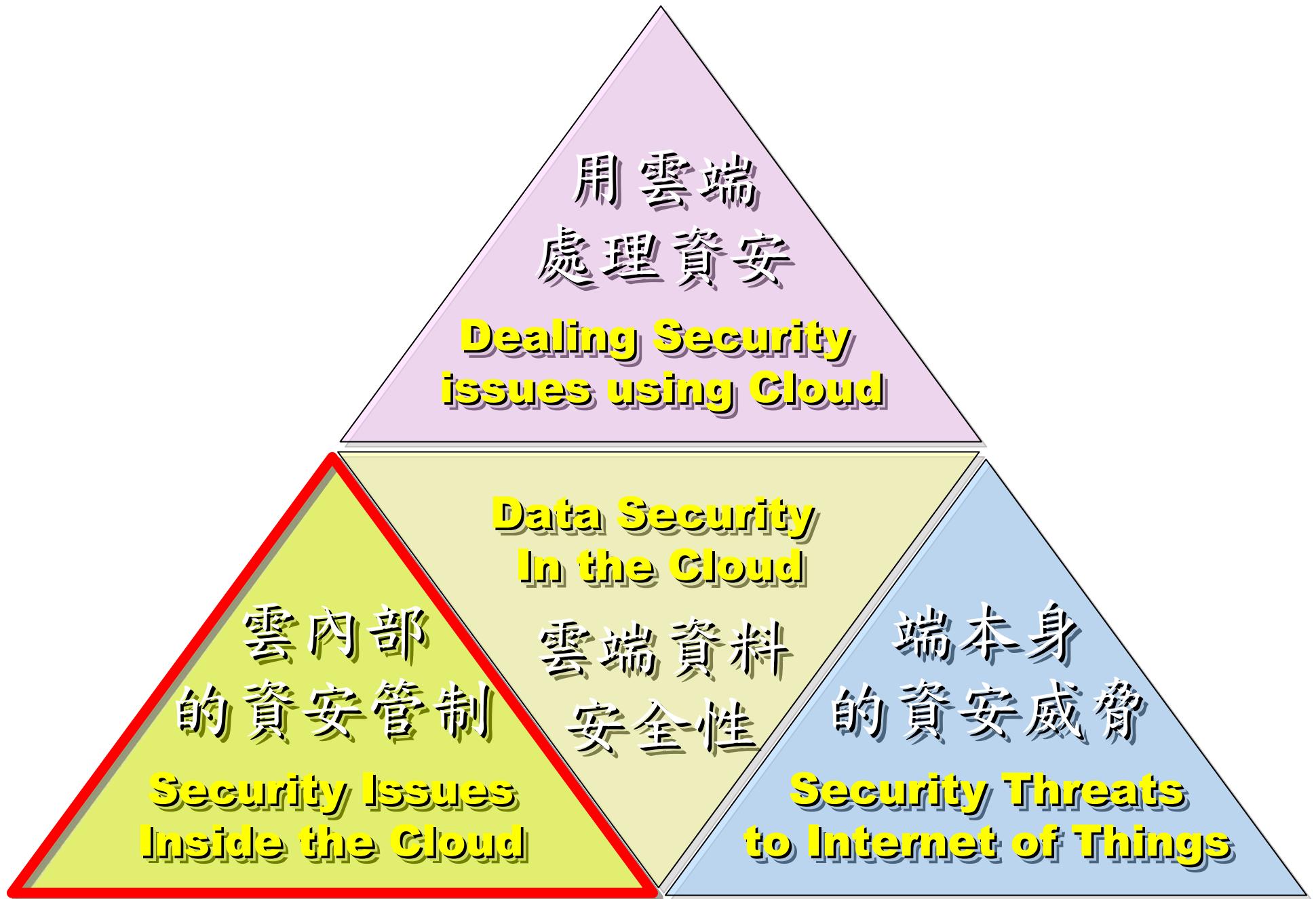
## U.S. Households with Smart Meters

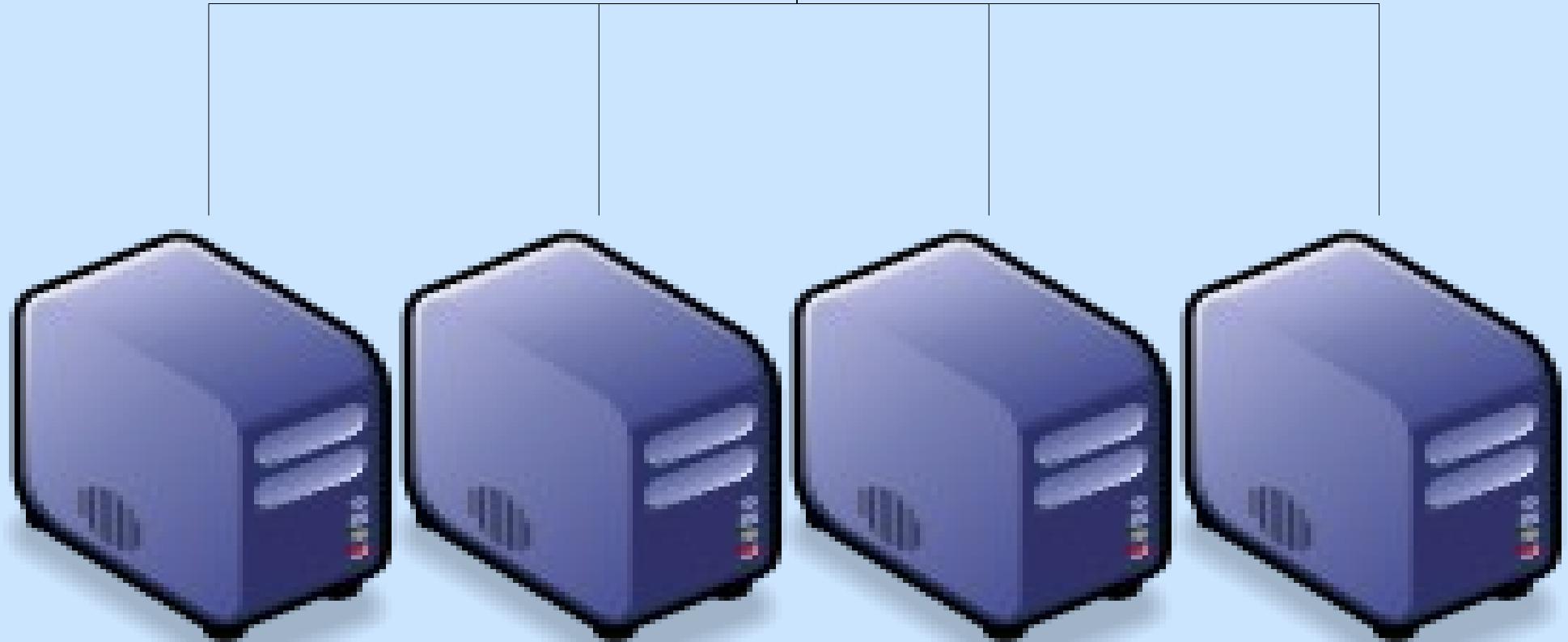


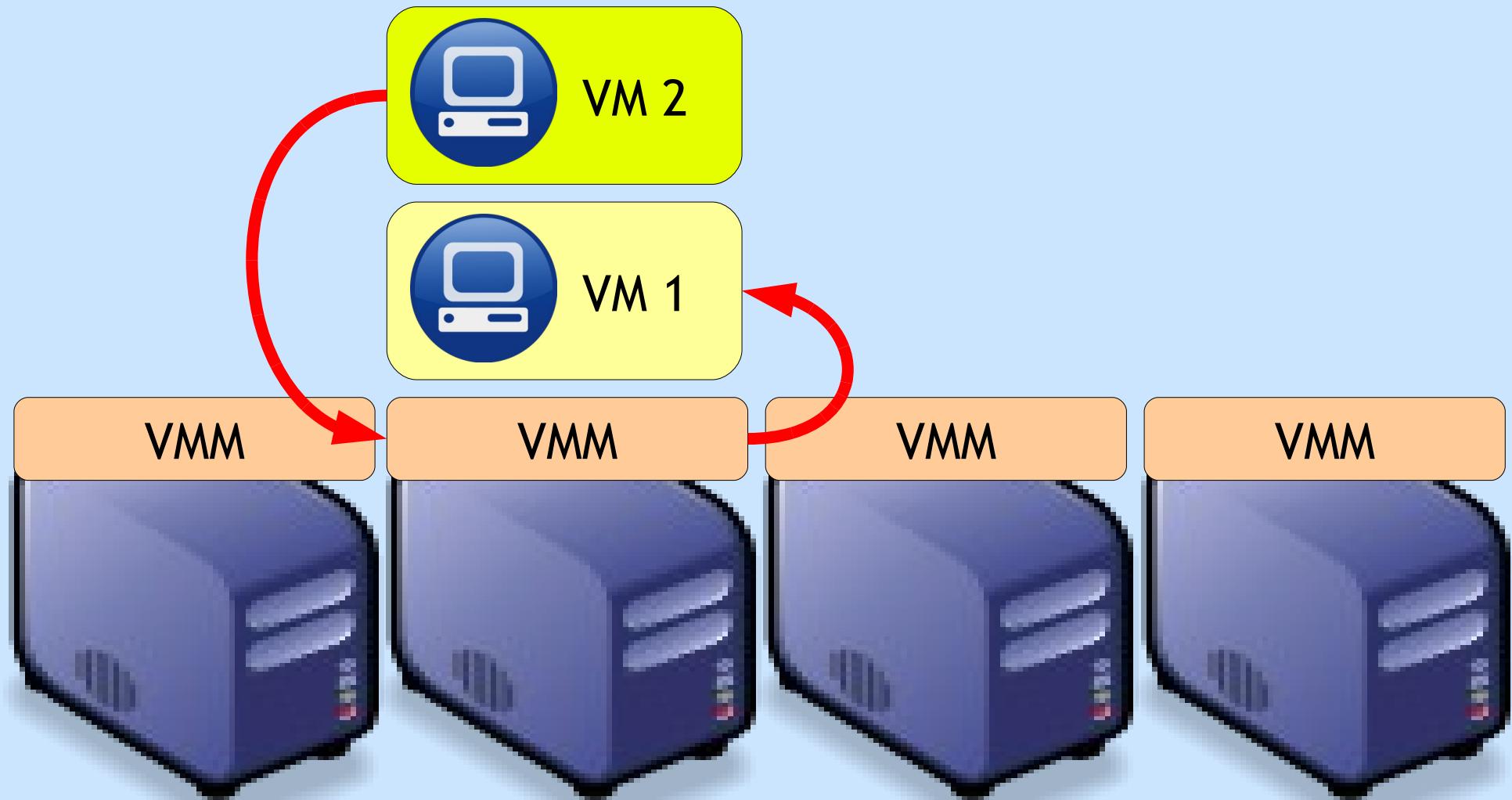
圖片來源：

[U.S. Unprepared for Internet Device Flood: Unaddressed Security Problems & Talent Drought Threaten Long-Term Commercial, Government Interests](#) 12  
By: Kurt Stammberger, CISSP, Adrian Turner and Mat Small, Mocana With: Rich Nass, Sarah Friar, Goldman Sachs

# 再來談談「雲的安全」

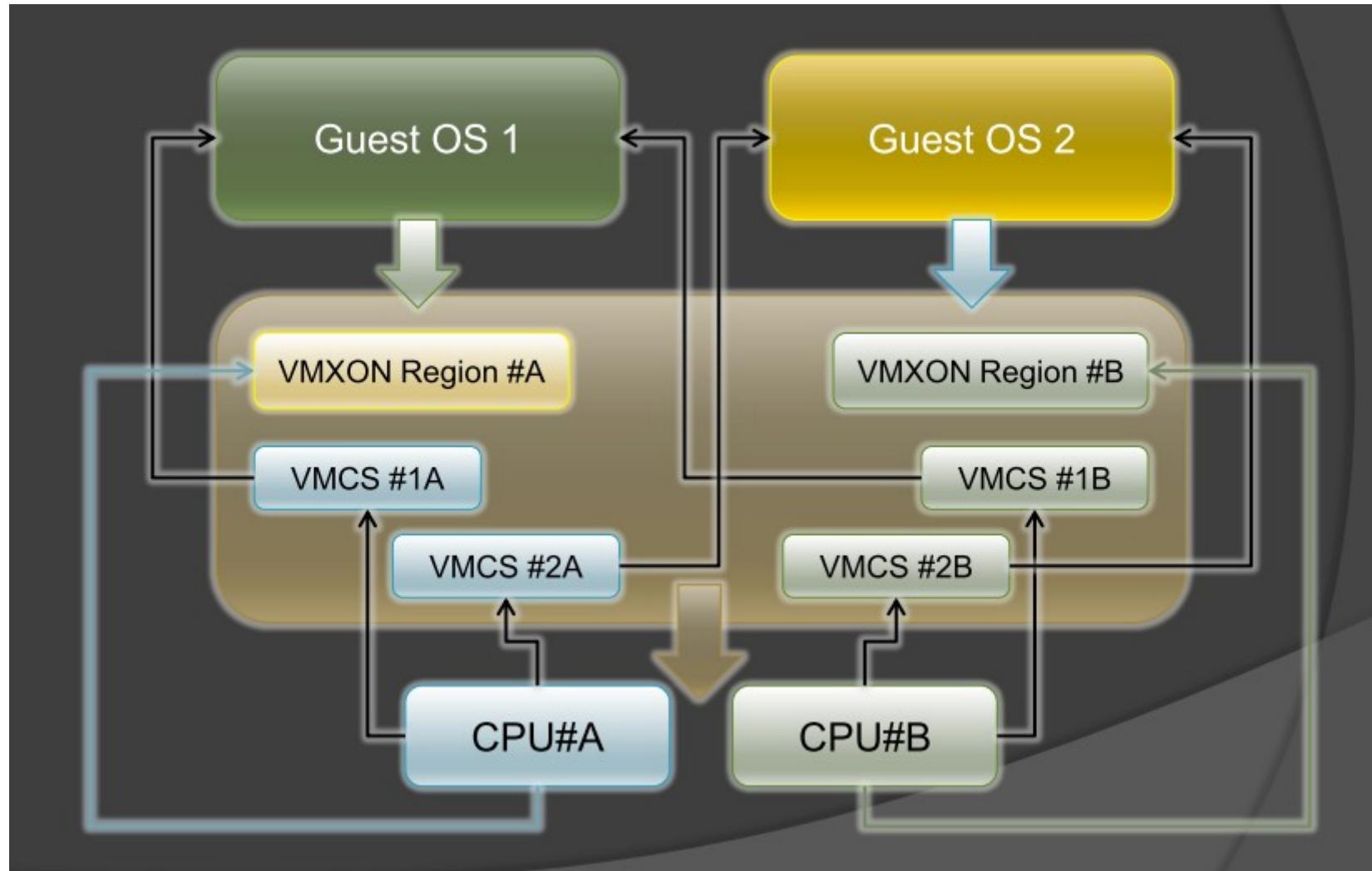






# 虛擬化衍生的新興資安問題

透過虛擬機器，竊取鍵盤輸入、植入後門.....



圖片來源：Hacks in Taiwan Conference 2010

[http://www.hitcon.org/hit2010/download/6\\_New%20Battlefield%20For%20Malware%20Game.pdf](http://www.hitcon.org/hit2010/download/6_New%20Battlefield%20For%20Malware%20Game.pdf)

王大寶 & PK / Hypervisor - New Battlefield For Malware Game 虛擬機 - 惡意程式攻防的新戰場 16



Virtual Switch

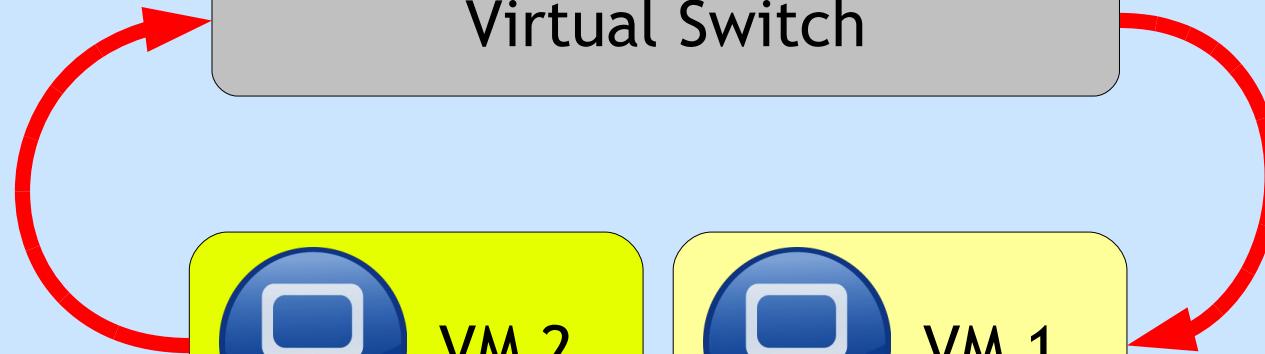


VMM

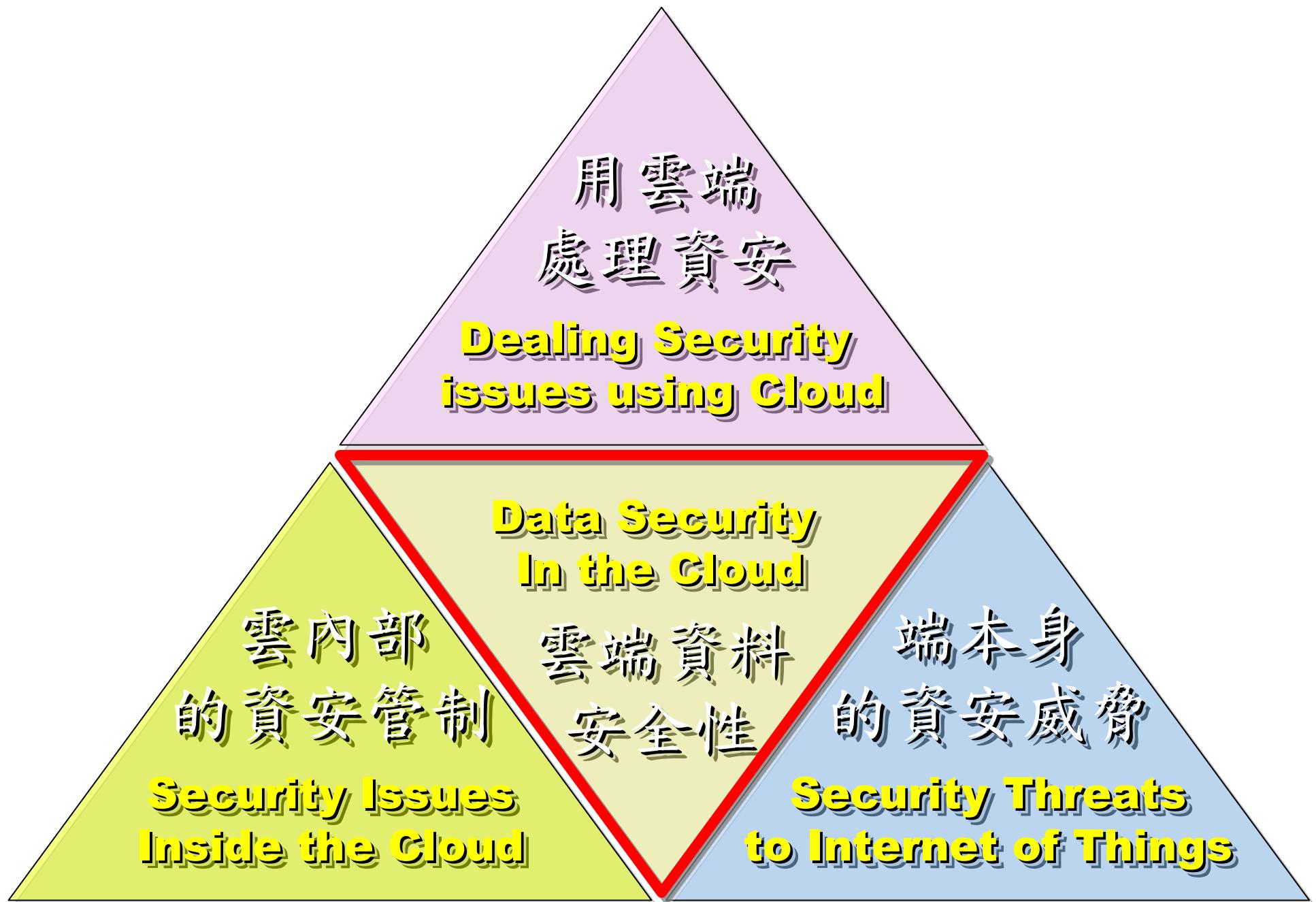
VMM

VMM

VMM



# 三談「資料安全」



# Ex. 無名照片外流、臉書個資外洩

轟動一時黑瀋會妹妹容瑄親密自拍照片外流



## WikiLeaks

“... could become as important a journalistic tool as the Freedom of Information Act.”

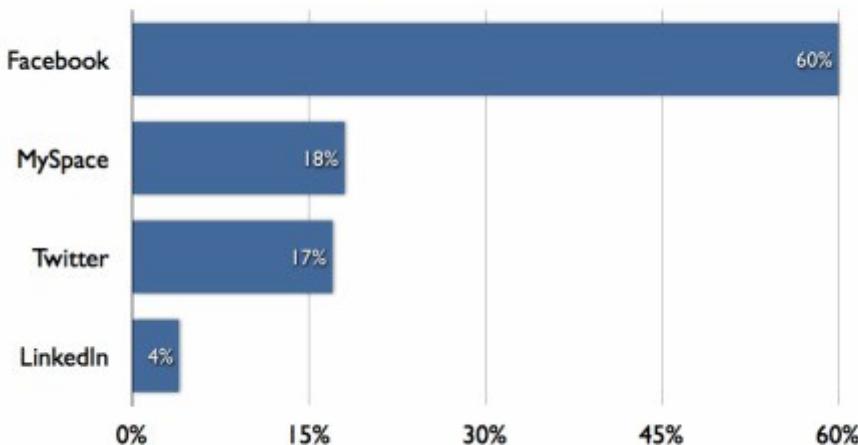
— Time Magazine

[Submit documents](#)

圖片來源：

[Wikileaks and Facebook Privacy / Security: Do we care?](#)

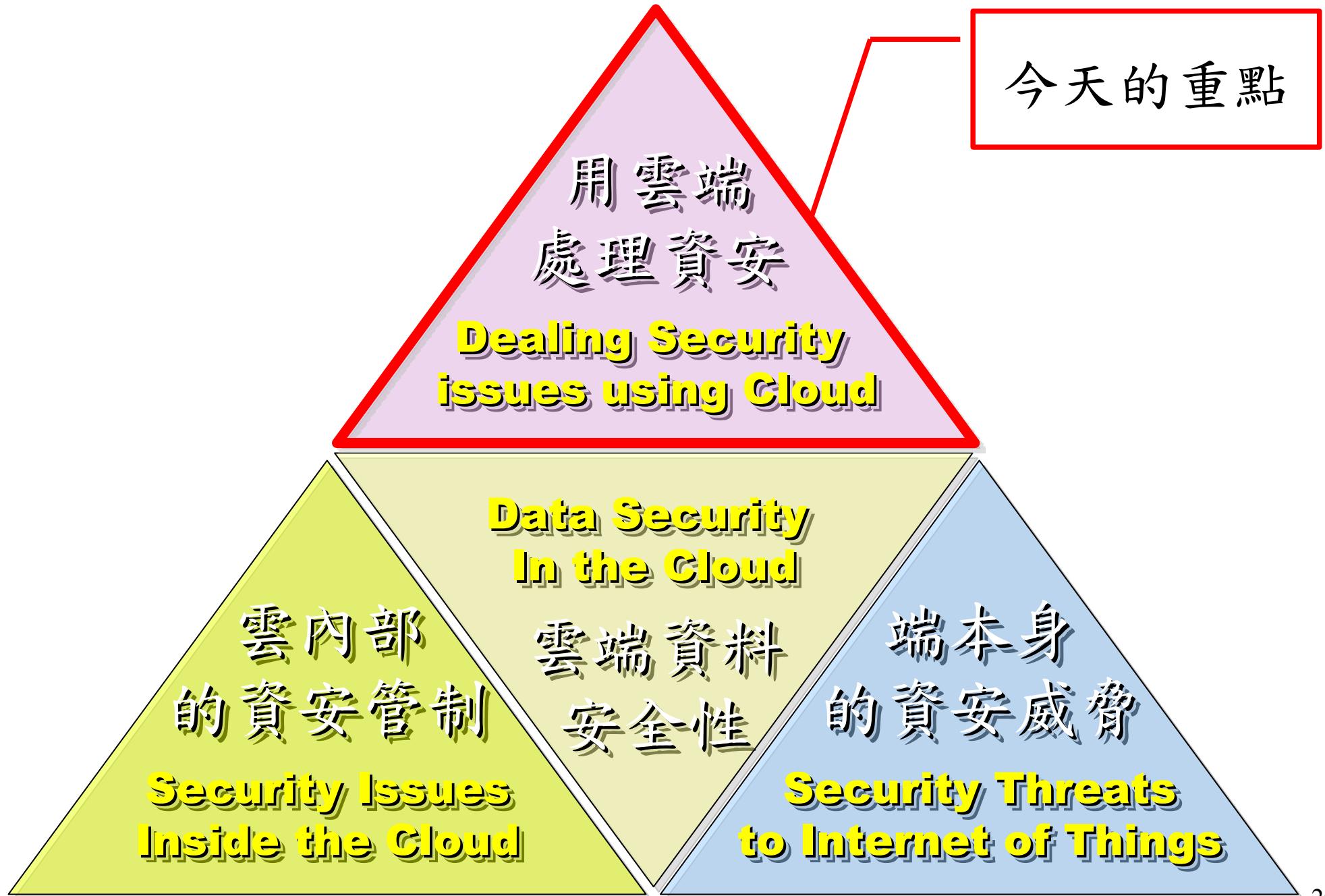
Which social network do you think poses the biggest risk to security?



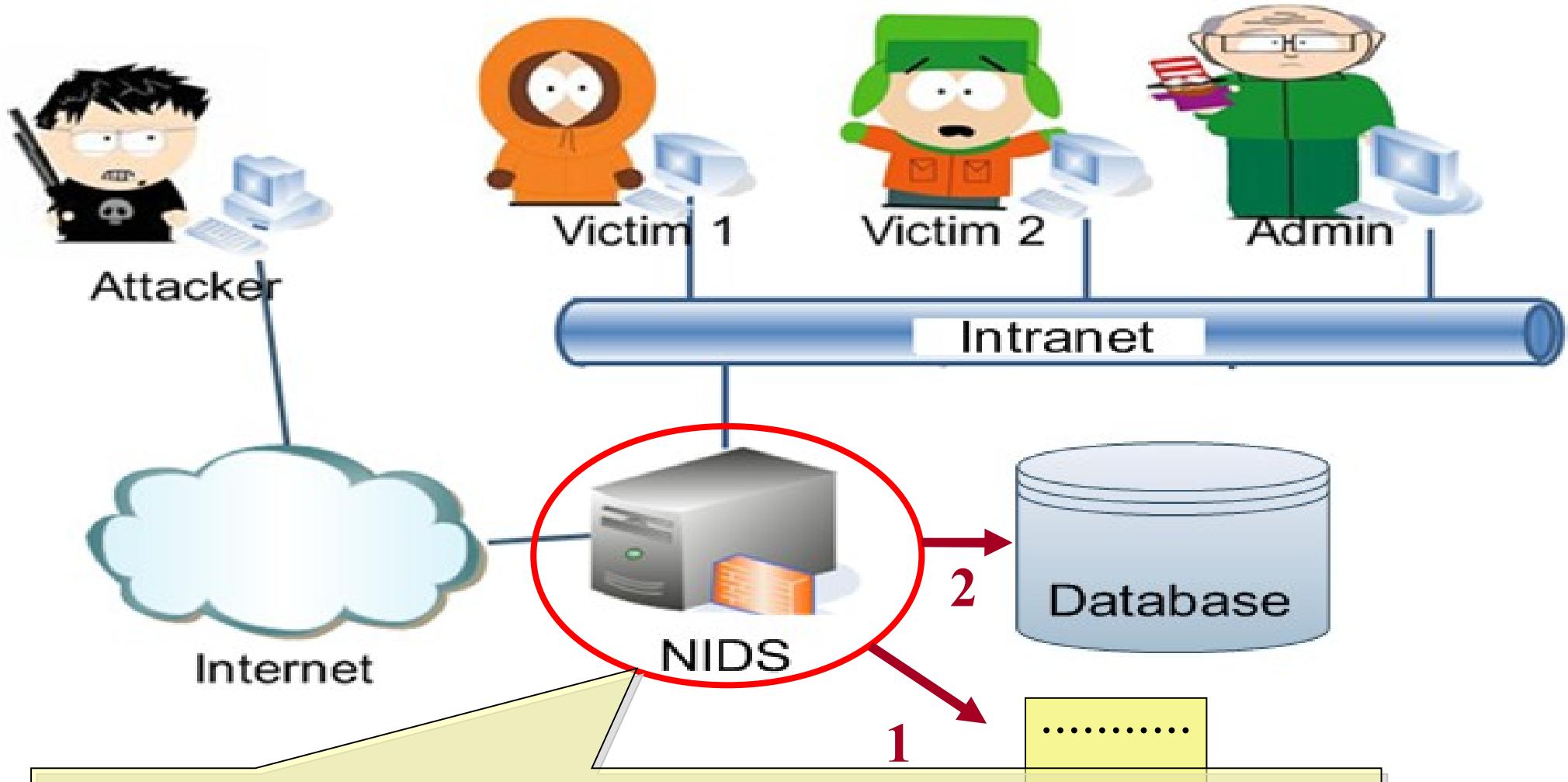
圖片來源：

[Report Ranks Facebook As Greatest Corporate Security Risk](#)  
<http://www.allfacebook.com/facebook-corporate-risk-2010-02>

# 進入今天的主題：用雲端處理傳統資安問題



# 使用入侵偵測系統 (NIDS) 來找出入侵訊息



當入侵偵測系統偵測到網路上有異常封包時，就會產生警訊以告知有攻擊發生。警訊通常有兩種形式：  
1. 紀錄成 log 檔 2. 紀錄到資料庫

# 傳統 NIDS 的警訊型態 (1) 紀錄在日誌檔內

## 入侵偵測系統所產生警訊日誌檔內一小段內容

[\*\*] [1:538:15] NETBIOS SMB IPC\$ unicode share access [\*\*]  
[Classification: Generic Protocol Command Decode] [Priority: 3]  
09/04/17:53:56.363811 168.150.177.165:1051 -> 168.150.177.166:139  
TCP TTL:128 TOS:0x0 ID:4000 IpLen:20 DgmLen:138 DF  
\*\*\*AP\*\*\* Seq: 0x2E589B8 Ack: 0x642D47F9 Win: 0x4241 TcpLen: 20

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.385573 168.150.177.164:1032 -> 239.255.255.250:1900  
UDP TTL:1 TOS:0x0 ID:80 IpLen:20 DgmLen:161  
Len: 133

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.386910 168.150.177.164:1032 -> 239.255.255.250:1900  
UDP TTL:1 TOS:0x0 ID:82 IpLen:20 DgmLen:161  
Len: 133

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.388244 168.150.177.164:1032 -> 239.255.255.250:1900  
UDP TTL:1 TOS:0x0 ID:84 IpLen:20 DgmLen:161  
Len: 133

[\*\*] [1:538:15] NETBIOS SMB IPC\$ unicode share access [\*\*]  
[Classification: Generic Protocol Command Decode] [Priority: 3]  
09/04/17:53:56.405923 168.150.177.164:1035 -> 168.150.177.166:139  
TCP TTL:128 TOS:0x0 ID:94 IpLen:20 DgmLen:138 DF  
\*\*\*AP\*\*\* Seq: 0x82073DFF Ack: 0x2468EB82 Win: 0x4241 TcpLen: 20

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.417045 168.150.177.164:45461 -> 168.150.177.1:1900  
UDP TTL:1 TOS:0x0 ID:105 IpLen:20 DgmLen:161  
Len: 133

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.420759 168.150.177.164:45461 -> 168.150.177.1:1900  
UDP TTL:1 TOS:0x0 ID:117 IpLen:20 DgmLen:160  
Len: 132

[\*\*] [1:1917:6] SCAN UPnP service discover attempt [\*\*]  
[Classification: Detection of a Network Scan] [Priority: 3]  
09/04/17:53:56.422095 168.150.177.164:45461 -> 168.150.177.1:1900  
UDP TTL:1 TOS:0x0 ID:118 IpLen:20 DgmLen:161  
Len: 133

[\*\*] [1:2351:10] NETBIOS DCERPC ISysTemActivator path overflow attempt little endian  
unicode [\*\*]  
[Classification: Attempted Administrator Privilege Gain] [Priority: 1]  
09/04/17:53:56.442445 198.8.16.1:10179 -> 168.150.177.164:135  
TCP TTL:105 TOS:0x0 ID:49809 IpLen:20 DgmLen:1420 DF  
\*\*\*A\*\*\*\* Seq: 0xF9589BBF Ack: 0x82CCF5B7 Win: 0xFFFF TcpLen: 20  
[Xref => http://www.microsoft.com/technet/security/bulletin/MS03-026.mspx][Xref =>  
http://cgi.nessus.org/plugins/dump.php3?id=11808][Xref => http://cve.mitre.org/cgi-  
bin/cvename.cgi?name=2003-0352][Xref => http://www.securityfocus.com/bid/8205]

[\*\*] [122:3:0] (portscan) TCP Portsweep [\*\*]  
[Priority: 3]  
09/04/17:53:56.499016 198.8.16.1 -> 168.150.177.166  
PROTO:255 TTL:0 TOS:0x0 ID:1750 IpLen:20 DgmLen:168

# 傳統 NIDS 的警訊型態 (2) 紀錄在資料庫內

以下為利用瀏覽器透過網頁方式呈現警訊資料庫的內容

Basic Analysis and Security Engine (BASE): Query Results - Mozilla

File Edit View Go Bookmarks Tools Window Help

Home | Search | AG Maintenance [ Back ]

Added 0 alert(s) to the Alert cache

Queried DB on : Thu October 14, 2004 22:04:44

Meta Criteria	any
IP Criteria	any
TCP Criteria	any
Payload Criteria	any

**Summary Statistics**

- Sensors
- Unique Alerts ( classifications )
- Unique addresses: source | destination
- Unique IP links
- Source Port: TCP | UDP
- Destination Port: TCP | UDP
- Time profile of alerts

Displaying alerts 1-50 of 81 total

ID	< Signature >	< Timestamp >	< Source Address >	< Dest. Address >	< Layer 4 Proto >
#0-(1-84)	[snort] NETBIOS SMB IPC\$ share unicode access	2004-10-08 11:25:41	192.168.1.100:1613	192.168.1.4:139	TCP
#1-(1-83)	[snort] NETBIOS SMB IPC\$ share unicode access	2004-10-08 11:25:31	192.168.1.100:1608	192.168.1.4:139	TCP
#2-(1-82)	[snort] NETBIOS SMB IPC\$ share unicode access	2004-10-08 11:25:05	192.168.1.100:1601	192.168.1.4:139	TCP
#3-(1-80)	[snort] (http_inspect) OVERSIZE CHUNK ENCODING	2004-10-04 22:25:41	192.168.1.4:42164	67.19.245.228:80	TCP
#4-(1-81)	[snort] (http_inspect) OVERSIZE CHUNK ENCODING	2004-10-04 22:25:41	192.168.1.4:42163	67.19.245.228:80	TCP

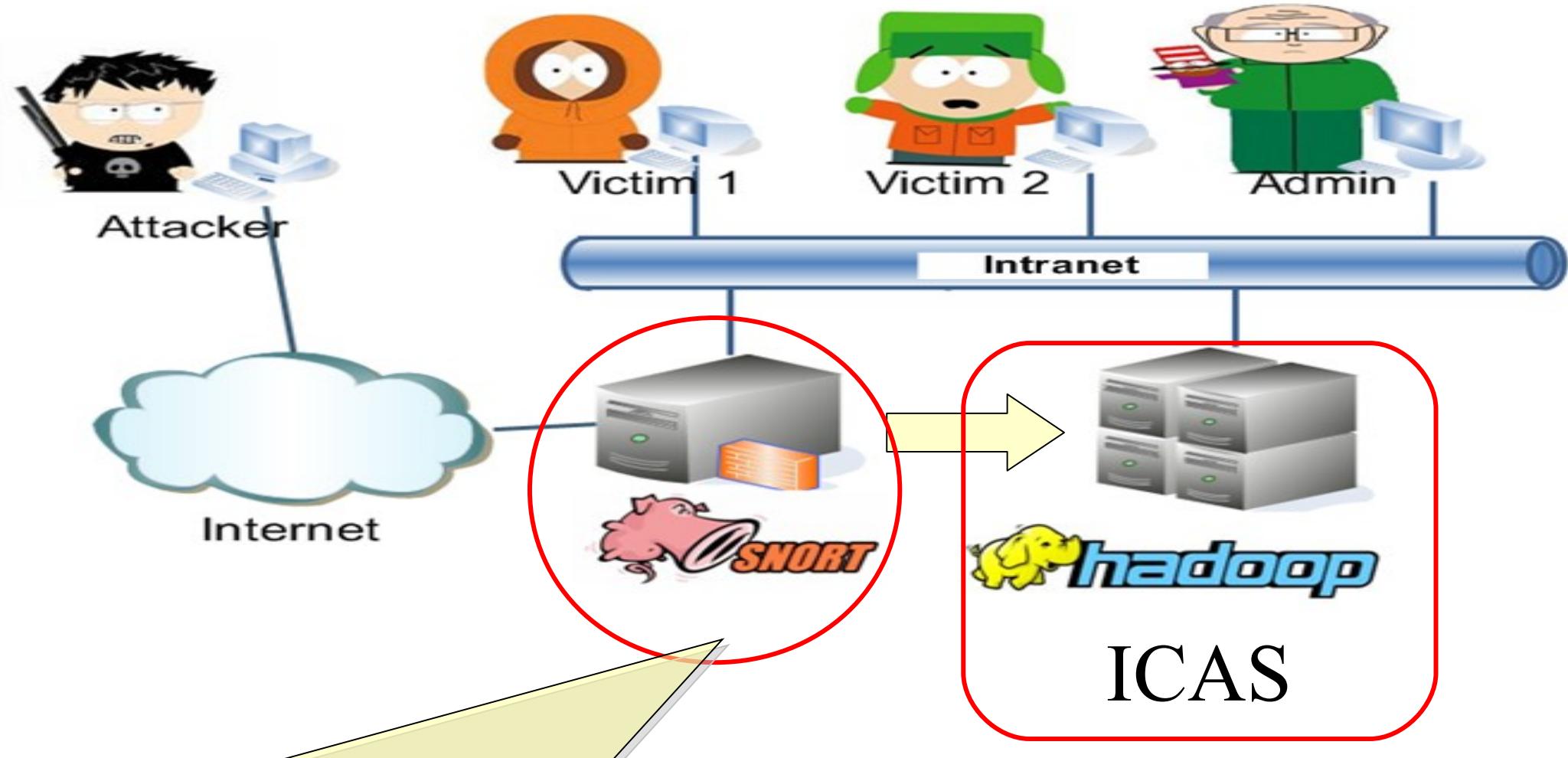
## 以上作法的缺點

- 警訊僅被『忠實』地被記錄下來，無法顯示彼此間的關聯性，因此系統管理者難以瞭解全部攻擊情形
- 過多的警訊，使得容易忽略重要內容
- 完全依賴單一台資料庫，當資料量一大，該台主機的讀寫效率將成為瓶頸

# 使用雲端運算的解決方案：ICAS

- ICAS, *IDS Cloud Analysis System*
- 利用雲端運算的特性提供以下好處
  - 對大量資料有高效率
  - 一般主機的叢集
  - 有錯誤容忍
- 分析演算法
  - 整合
  - 關聯

# 透過 ICAS 協助分析 IDS 的警訊

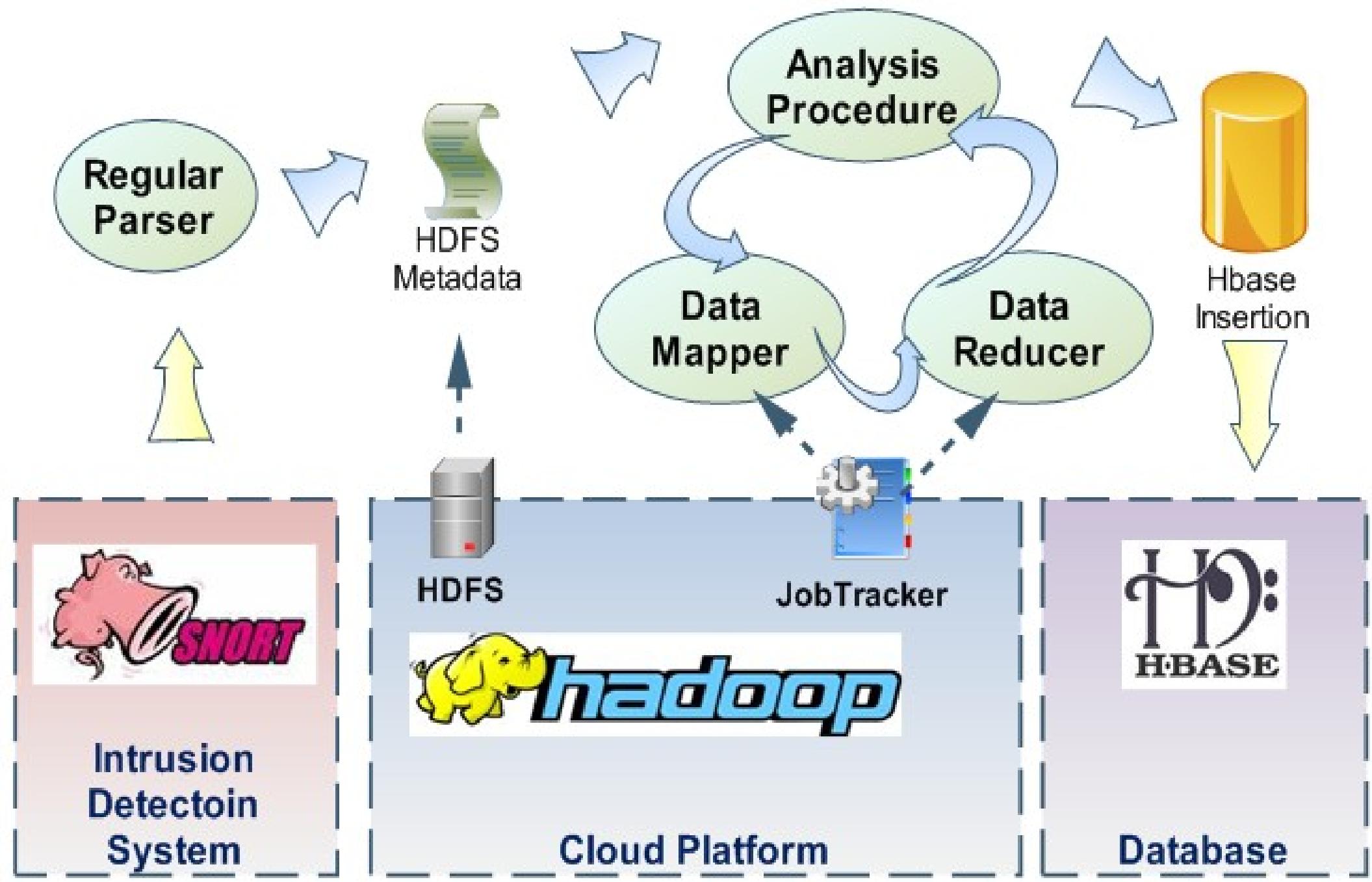


可多個 NIDS 共同產生警訊後，傳送至 ICAS ，分析演算法  
目前有 ICAS-I 及 ICAS-II

# ICAS-I

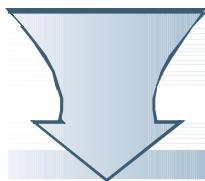
- 將任意個原始警訊檔上傳到運行 ICAS-I 演算法的 Hadoop 檔案系統空間（HDFS）
- 利用 Hadoop 的 MapReduce 平台架構所設計的演算法來分析資料
- 分析完後的資料塞入分散式資料庫 HBase 內

# ICAS-I 流程圖



# ICAS-I 整合後的警訊結果

Destination IP	Attack Signature	Source IP	Destination Port	Source Port	Packet Protocol	Timestamp
Host_1	Trojan	Sip1	80	4077	tcp	T1
Host_1	Trojan	Sip2	80	4077	tcp	T2
Host_1	Trojan	Sip1	443	5002	tcp	T3
Host_2	Trojan	Sip1	443	5002	tcp	T4
Host_3	D.D.O.S	Sip3	53	6007	udp	T5
Host_3	D.D.O.S	Sip4	53	6008	tcp	T5
Host_3	D.D.O.S	Sip5	53	6007	udp	T5
Destination IP	Attack Signature	Source IP	Destination Port	Source Port	Packet Protocol	Timestamp



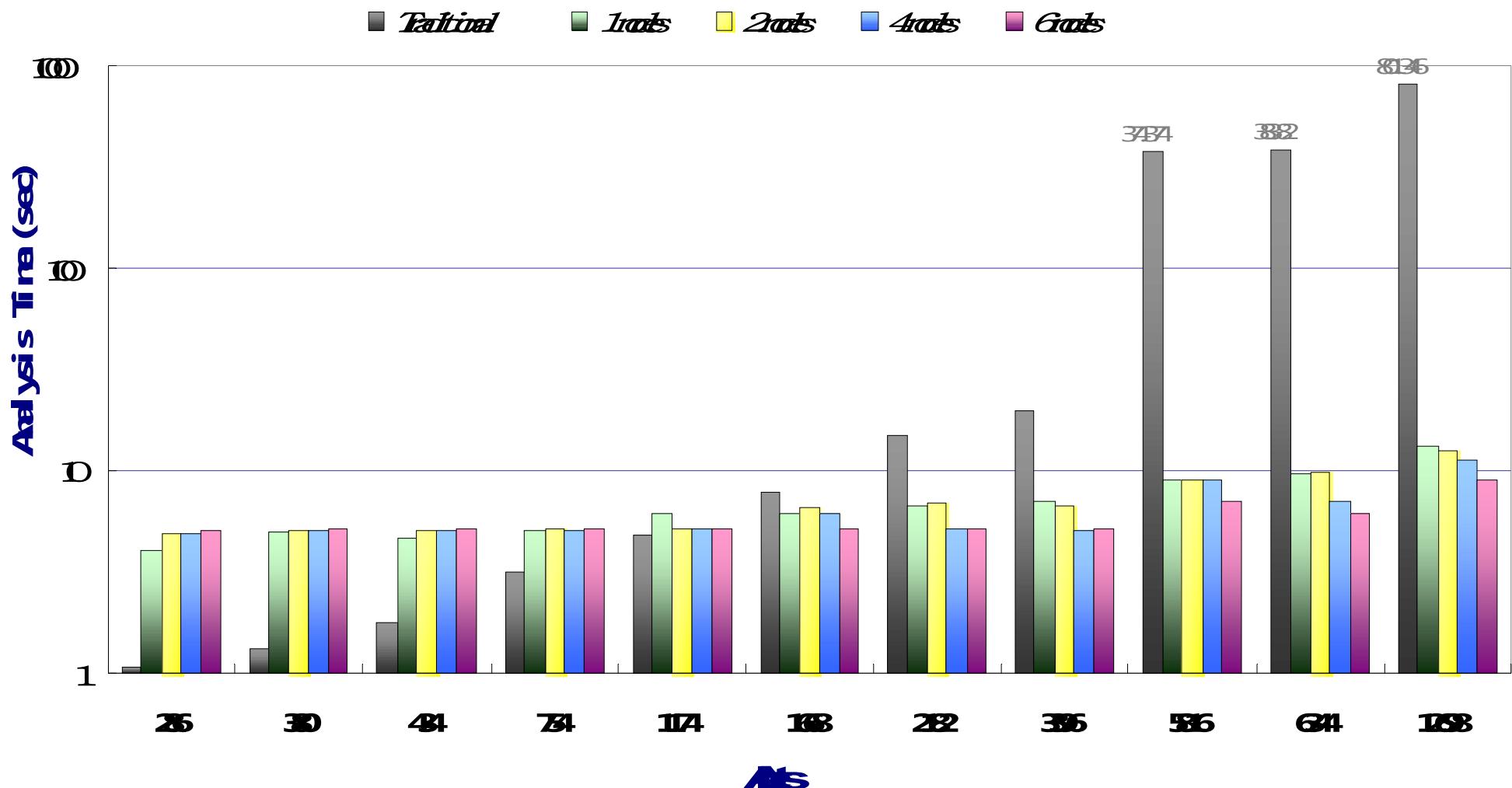
Key	Values
Host_1	Trojan
Host_1	Sip1,Sip2
Host_2	Trojan
Host_2	Sip1
Key	Values

# ICAS-I 效能數據的環境

- Machine:
  - CPU : Intel quad-core, Memory : 2 GB,
- OS : Linux : Ubuntu 8.04 server
- Software : version
  - Hadoop : 0.16.4
  - Hbase : 0.1.3
  - Java : 6
- Alerts Data Sets
  - MIT Lincoln Laboratory, Lincoln Lab Data Sets
  - Computer Security group at UCDavis, tcpdump file

# ICAS-I 效能分析時間圖

## The Consuming Time of Each Number of Data Sets



# ICAS-I 效能數據表

## Throughput Data Overall

Original Alerts	Analysis Time (sec)					Results	Reduction Rate
	Traditional	1 nodes	2 nodes	4 nodes	6 nodes		
286	1.068	4.087	4.869	4.864	5.077	30	89.51%
380	1.333	4.94	5.069	5.067	5.097	11	97.11%
434	1.76	4.61	5.066	5.068	5.09	9	97.93%
754	3.145	5.066	5.079	5.038	5.096	16	97.88%
1174	4.73	6.066	5.093	5.089	5.097	33	97.19%
1668	7.909	6.07	6.56	6.071	5.082	16	99.04%
2182	14.949	6.671	6.95	5.166	5.088	16	99.27%
3396	19.901	7.053	6.654	5.076	5.091	68	98.00%
5816	374.374	9.081	9.076	9.07	7.076	66	98.87%
6344	383.82	9.68	9.872	7.069	6.069	72	98.87%
12698	801.346	13.096	12.367	11.367	9.083	36	99.72%

# ICAS-II

- ICAS-I 僅將資料塞入資料庫，然而還是文字的敘述
- ICAS-II 將輸入的任意多個警訊整合成一張警訊關聯圖
- 資料的來源可以透過以下兩種方式上傳到分析平台
  - 系統自動設定以 SCP 傳送到 ICAS 工作目錄
  - 管理者透過 ICAS 網頁上傳

透過網頁上傳的方式將 snort 的警訊檔送至 icas 分析

The screenshot shows a web browser window with the URL <http://secuse.nchc.org.tw/icas/checkpoint.php>. The page title is 'ICAS'. A message from Google Translate asks if it should translate the page. The main content area has a yellow background and contains the following text:

一旦選定需分析的日誌檔  
後，按下 upload 鈕，系  
統進行上傳→分析→繪圖  
等步驟

Please upload your Snort Log files  
(ex: /var/log/snort/alert...).

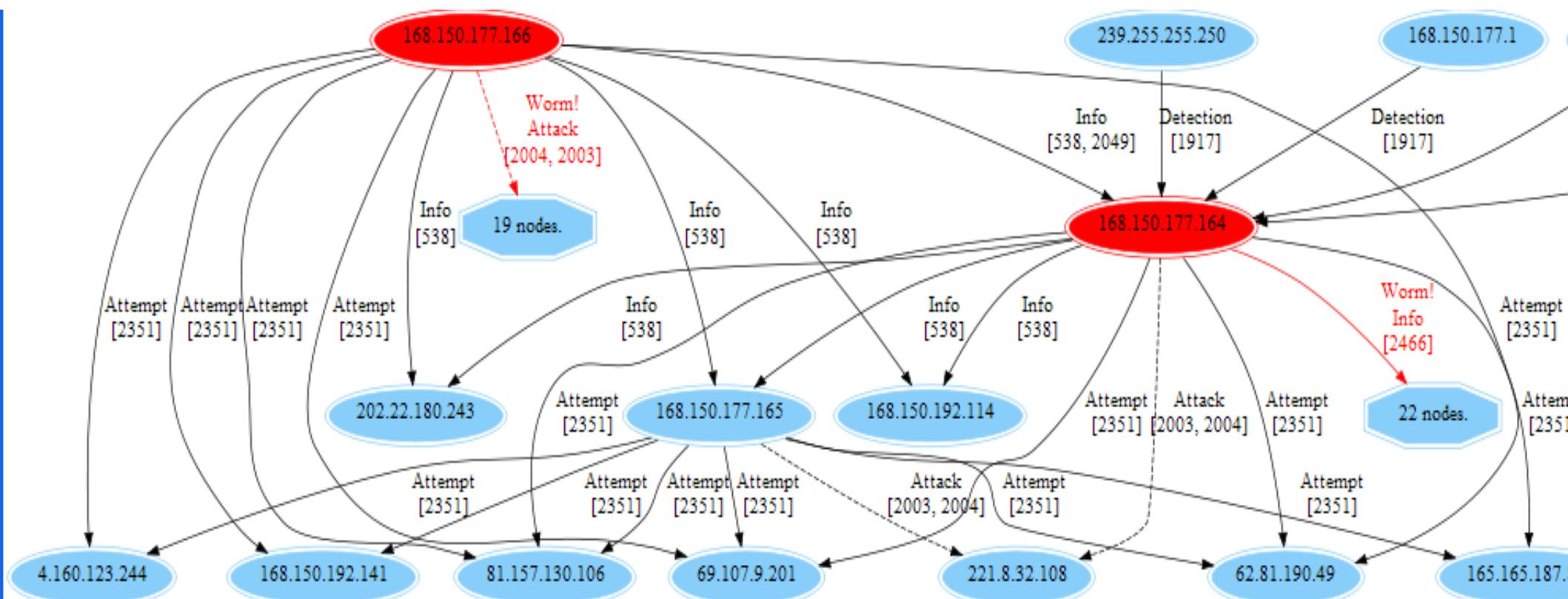
透過網頁上傳的方式將 snort 的警訊檔送至 icas 分析

瀏覽... upload

go to final report directory

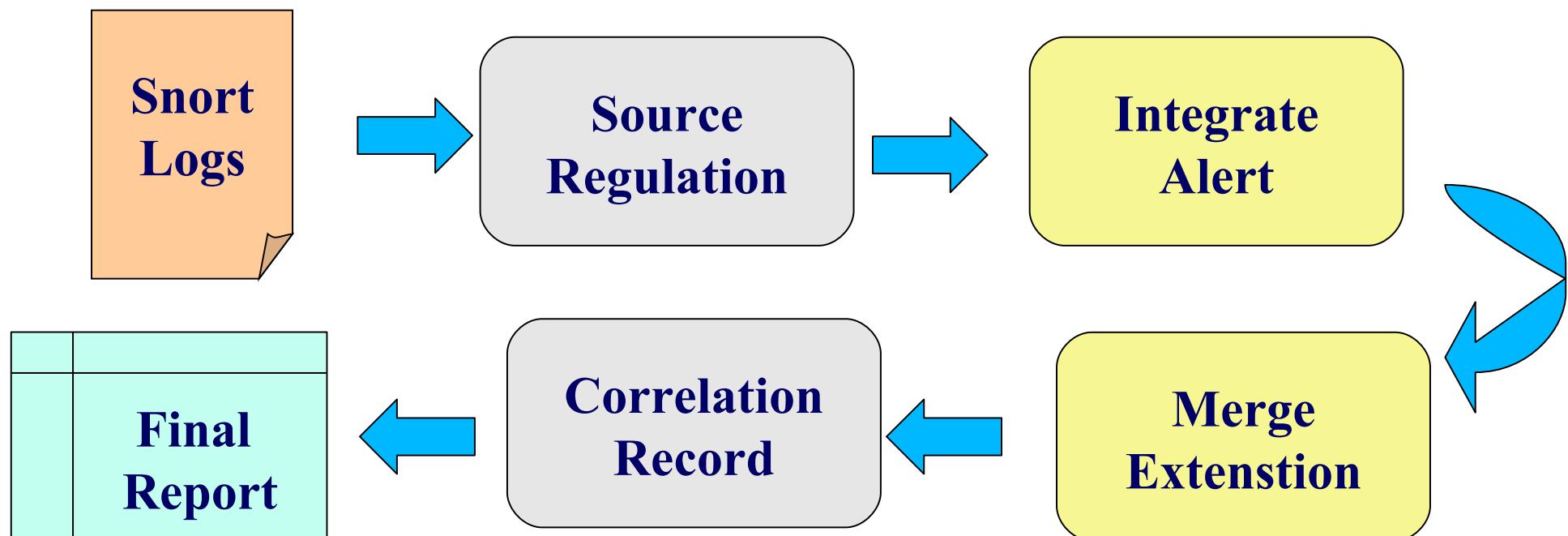
# ICAS-II 所產生的報表：警訊關聯圖

- 經過 ICAS-II 分析後，可以得到此警訊關聯圖。
- 圖中橢圓形代表節點，箭頭及線上文字代表攻擊方向與攻擊方法。
- 標為紅色則是經過系統分析之後，被判定有攻擊行為的節點與方法。
- 此圖說明 IP 168.150.177.166 與 168.150.177.164 有進行蠕蟲的攻擊行為



# ICAS-II 的分析流程

- Hadoop v 0.20



# ICAS-II 結論

- ICAS-II 可經過警訊的來源、目的、攻擊事件綜合分析
  - 提供巨觀攻擊關聯圖來瞭解攻擊事件的始末
  - 自動透過標記顏色的方法將較高危險的事件呈現出來。
- ICAS-II 尚在整合關聯式資料庫，因此還未進行數據量測

# ICAS 總結

- 雲端運算處理資料格式相似且資料量大的情況下，能展現其效益
- 提供高容錯率、低獨占系統資源、多工作同時執行等能力
- 可搭配其他軟體作即時的警訊資料呈現， ICAS 可補充分析後資料的部份
- 未來工作
  - 整合多種資料來源平台
  - 產生更詳細與人性化的分析資料



## Questions?

Slides - <http://trac.nchc.org.tw/cloud>

**Jazz Wang  
Yao-Tsung Wang  
jazz@nchc.org.tw**

